

12 Deep Ensemble Learning for High-dimensional
13 Subsurface Fluid Flow Modeling

14 Abouzar Choubineh^{a,b}, Jie Chen^b, David A. Wood^c, Frans Coenen^a, Fei
15 Ma^b

^a*Department of Computer Science, University of Liverpool, Liverpool, United Kingdom*

^b*Department of Applied Mathematics, Xi'an Jiaotong-Liverpool University, Suzhou, China*

^c*DWA Energy Limited,, Lincoln, United Kingdom*

16 **Abstract**

The accuracy of Deep Learning (DL) algorithms can be improved by combining several deep learners into an ensemble. This avoids the continuous endeavor required to adjust the architecture of individual networks or the nature of the propagation. This study investigates prediction improvements possible using Deep Ensemble Learning (DEL) to determine four distinct multiscale basis functions in the mixed Generalized Multiscale Finite Element Method (GMsFEM), involving the permeability field as the only input. 376,250 samples were initially generated, filtered down to 367,811 after data pre-processing. A standard Convolutional Neural Network (CNN) named SkiplessCNN and three skip connection-based CNNs named FirstSkipCNN, MidSkipCNN, and DualSkipCNN were developed for the base learners. For each basis function, these four CNNs were combined into an ensemble model using linear regression and ridge regression, separately, as part of the stacking technique. A comparison of

the coefficient of determination (R^2) and Mean Squared Error (MSE) confirms the effectiveness of all three skip connections in enhancing the performance of the standard CNN, with DualSkip being the most effective among them. Additionally, as evaluated on the testing subset, the combined models meaningfully outperform the individual models for all basis functions. The case that applies linear regression delivers R^2 ranging from 0.8456 to 0.9191 and MSE ranging from 0.0092 to 0.0369. The ridge regression case achieves marginally better predictions with R^2 ranging from 0.8539 to 0.922, and MSE ranging from 0.009 to 0.0349 because its solution involves more evenly distributed weights.

17 *Keywords:*

18 Deep ensemble learning, Big data, Linear/ridge regression,
19 Convolutional neural network, Skip connection, Subsurface fluid flow,
20 mixed generalized multiscale finite element method.

21 **1. Introduction**

22 The wide range of Machine Learning (ML) algorithms available all
23 have to contend with reducible and irreducible errors. The latter is typi-
24 cally a consequence of noise within the datasets being evaluated and can-
25 not be addressed by the ML models themselves. On the other hand, bias
26 and variance combine to generate reducible errors, which can be effec-
27 tively reduced by the algorithmic actions of the ML. Bias errors are a con-
28 sequence of the differences between predicted and actual dependent vari-

29 able values generated with a training subset of samples. Variance errors
 30 result from small fluctuations in the training subsets actual values. Math-
 31 ematically, we may assume that there is an input vector \mathbf{X} (here the perme-
 32 ability field) that influences an output vector \mathbf{Y} (here the basis function).
 33 The function $f(\mathbf{X})$ denotes the correct relationship between the input and
 34 output, but it is accompanied by some noise that can be represented by σ_{ϵ}^2
 35 that constitutes the irreducible error:

$$\mathbf{Y} = f(\mathbf{X}) + \sigma_{\epsilon}^2 \quad (1)$$

36 ML models strive to determine the best function $\hat{f}(\mathbf{X})$ that can predict
 37 the true underlying function $f(\mathbf{X})$ as precisely as possible. Given the Total
 38 Error (TE) as $TE = E[(\mathbf{Y} - \hat{f}(\mathbf{X}))^2]$:

$$TE = [E\hat{f}(\mathbf{X}) - f(\mathbf{X})]^2 + E[\hat{f}(\mathbf{X}) - E\hat{f}(\mathbf{X})]^2 + \sigma_{\epsilon}^2 \quad (2)$$

$$TE = bias^2 + variance + irreducible error \quad (3)$$

39 Simpler models tend to generate high bias accompanied by low vari-
 40 ance. On the other hand, more elaborate models tend to generate lower
 41 bias accompanied by higher variance. Linear regression, for example,
 42 has a high bias since it tends to oversimplify and, therefore, cannot accu-
 43 rately capture the relationship between input variables and output data.
 44 In contrast, artificial neural networks involving multiple hidden layers

45 and many nodes tend to generate substantial variance, because they tend
46 to overfit training datasets, making it difficult for the trained models to
47 be generalized and accurately predict unseen data. High bias is typically
48 a consequence of models underfitting a dataset, whereas high variance is
49 typically a consequence of models overfitting a dataset. As a modeling
50 strategy, it makes sense, therefore, to attempt to trade off bias and vari-
51 ance errors to assist the trained models in being applied in a more gener-
52 alized way, thereby more accurately predicting data not seen during the
53 training/validation process.

54 Ensemble learning, where by a number of base learners are combined
55 into an “ensemble” to produce a single model whose predictive or clas-
56 sification accuracy is better than that of the individual component base
57 learners, is a well-established technology (Wang et al., 2011, 2014; Nguyen
58 and Logofătu, 2018; Kanda et al., 2020; Verma and Chandra, 2023; Sahin
59 and Demir, 2023). The majority of published research directed towards
60 ensemble learning has been founded on traditional ML techniques; for
61 example *Random Forests* (Ho, 1998). Such established mechanisms repre-
62 sent Shallow Ensemble Learning (SEL). The alternative is Deep Ensemble
63 Learning (DEL), which combines a number of deep learners into an en-
64 semble. Although Deep Learning (DL) algorithms tend to generate fewer
65 prediction errors than ML methods when applied to many datasets, there
66 is scope to further improve their accuracy. Combining several deep learn-
67 ers into an ensemble is one way to potentially achieve this. Moreover, in

68 a set of deep models, the different strengths of each DL model may com-
69 plement one another, and weaknesses cancel each other out. Unlike SEL,
70 DEL has received little attention to date. Rather, attempts to improve DL
71 accuracy have been focused on optimizing the various control-parameter
72 values used by deep learners, for example by adjusting the architecture of
73 a DL network or the nature of the propagation applied. This study chal-
74 lenges this view by investigating the performance of DEL. It is not possi-
75 ble to use common ML techniques such as Fully Connected (FC) models
76 for the problem investigated in this paper, which is mapping an input of
77 100×9 to an output of 900×1 . This is mainly because the input is a
78 2D tensor. Among DL techniques, recurrent neural networks are usually
79 applied to video, sound, or text data. On the other hand, Convolutional
80 Neural Networks (CNNs) are specifically designed for problems with 2D
81 arrays. Therefore, base learners are selected from complex CNN models
82 with different variants. These base learners are then combined using two
83 regression models (linear regression and ridge regression), separately.

84 A key motivation for this study is to develop DEL models to assist in
85 the prediction of fluid-flow characteristics in subsurface reservoirs. This
86 is of interest to provide more detailed insights to the flow of fluids into
87 producing wells penetrating oil and gas reservoirs, the seepage of fluids
88 through soil, and land subsidence as a consequence of groundwater and
89 oil and gas extraction. With respect to oil and gas reservoirs, the main
90 goal is to predict the performance of reservoirs at any future point in time

91 and to optimize petroleum fluid recovery under different operating condi-
92 tions. Fluid flow in petroleum reservoirs is typically modelled using a set
93 of non-linear Partial Differential Equations (PDEs). In general, these equa-
94 tions can be solved analytically (exact solutions) or numerically (approx-
95 imate solutions). However, for reservoir simulation models comprised of
96 many thousands of grid cells the datasets generated to do this involve hun-
97 dreds of thousands of data points leading to long and tedious calculations.
98 Moreover, the datasets involved are typically too large for ML models to
99 handle. This study provides an innovative solution by developing efficient
100 and reliable DEL models to assist with a specific time-consuming aspect
101 of fluid-flow simulation.

102 The developed DEL models are applied to a large dataset associated
103 with subsurface fluid flow modeling. This dataset can be modeled in dif-
104 ferent ways but a mixed Generalized Multiscale Finite Element Method
105 (GMsFEM) is used to generate multiscale data formats. 376,250 samples
106 (data records) were generated with a 2D permeability field (in a Cartesian
107 coordinate system) representing the input variable and the multiscale ba-
108 sis functions as the output. Traditionally, a substantial number of PDEs
109 need to be solved to produce these functions, but this involves consid-
110 erable time and computational effort. Deep ensemble learning made up
111 of several deep base models offers an innovative and more effective alter-
112 native to PDE solvers, specifically with respect to determining accurate
113 reservoir pressure distributions to improve estimates of resource recovery

114 factors from oil/gas reservoirs.

115 Following on from this introduction, Section 2 considers ensemble mod-
116 eling along with published research relevant to it. A review of the mixed
117 GMsFEM is given in Section 3. The characteristics of stacking CNN en-
118 semble models are presented in Section 4. Prediction error evaluation of
119 applications of the developed DEL models to the subsurface fluid flow
120 dataset, in terms of the coefficient of determination (R^2) and Mean Squared
121 Error (MSE), is presented in Section 5. Section 6 discusses the implications
122 of the results generated leading to the conclusions drawn (Section 7).

123 2. Ensemble Modeling with Its Related Research

124 Ensemble systems can be categorised according to the method by which
125 the ensemble learning is achieved:

- 126 1. **Boosting:** Boosting (Schapire, 1990; Kumari and Toshniwal, 2021) is
127 a sequential method. The different base learners are dependent on
128 each other. The aim is to fit models in steps such that model training
129 at each step is influenced by the models constructed in the previous
130 steps. Each step is focused on examples in the dataset that have been
131 poorly predicted by the previous steps.
- 132 2. **Bagging:** Bagging (standing for bootstrap aggregating) (Breiman,
133 1996a; Tüysüzöğlü and Birant, 2020) is a parallel approach to en-
134 semble learning where multiple models are generated using the same
135 ML algorithm but with different portions of the training data, which

136 are then merged to produce a single more robust model than the in-
137 dividual base models. Multiple training subsets (bootstrap samples)
138 are randomly selected from the initial training dataset with replace-
139 ment (a single row of the initial data might be chosen zero, one, two,
140 or even more times). Each model is developed from one subset, re-
141 sulting in an ensemble of several models. The final prediction is
142 obtained by averaging (or simply ranking) all the predictions of the
143 different learners.

144 3. **Stacking:** Unlike boosting and bagging, stacking (Breiman, 1996b;
145 Yin et al., 2021) uses base learners generated by different machine
146 learners. The voting ensemble represents a simple stacking method
147 in which a statistical mechanism is used to combine different types of
148 ML models, such as decision trees and support vector machines. No
149 matter how well the individual ML models perform on the training
150 dataset, they all contribute equally to the merged model. One can
151 consider the simple average of the predictions from the underlying
152 ML models. However, using a weighted average ensemble makes the
153 results more sensitive to the prediction errors generated by each con-
154 tributing ML model. A further improvement can be made through
155 stacked generalization, which applies a ML model to learn how to
156 best combine the predictions derived from the base learners. It does
157 this by first developing base models using the training dataset in-
158 puts. It then feeds the underlying ML models into a meta-learner,

159 which attempts to make a new model using the predictions of the
160 weak learners based on new data.

161 Ensemble techniques are now widely applied in various engineering
162 and geoscience disciplines. For example, three models of Bayesian, func-
163 tional, and meta-ensemble were applied to Land Subsidence Susceptibil-
164 ity (LSS) mapping (Oh et al., 2019). The models split the dataset 50:50
165 between training and testing subsets with errors measured in relation to
166 operating-characteristic curve. The ensemble including the logit boost
167 model delivered the most accurate (91.44%) LSS maps.

168 Slope stability predictions were generated using a hybrid stacking en-
169 semble method (Kardani et al., 2021). An artificial bee colony optimizer
170 was applied to identify the optimal combination of base classifiers (ensem-
171 ble level 0). These were then used to develop an effective meta-classifier
172 (ensemble level 1), considering eleven separate tuned ML models. Fi-
173 nite element analysis was employed to create a synthetic database (150
174 records) for training the models. The trained models were then applied
175 to predict 107 naturally occurring slope cases to test model performances.
176 The hybrid-stacking ensemble model generated less errors than each ML
177 model used in isolation.

178 Ensemble random forest, ensemble Gradient Boosted Regression Tree
179 (GBRT) and multiLayer perceptron neural network were applied to model
180 the spatial extent of landslides in Norway (Liu et al., 2021). Eleven landslide-
181 influencing factors were considered related to geomorphologic, geologic,

182 geo-environmental, and anthropogenic effects. 3,399 positive landslide
183 records and 6,798 non-landslide were considered. Seventy percent of the
184 data records in each of these two categories were selected for training the
185 models. The remaining thirty percent of the data records were used to
186 test the trained models. Slope angle was confirmed by the models to be
187 the most important influencing factor. The ensemble GBRT model outper-
188 formed the other ensemble models, achieving a 95% probability of detect-
189 ing landslides in that region.

190 Support vector machine, multilayer perceptron, random forest, Adap-
191 tive Boosting (AdaBoost), and extreme gradient boosting were used to de-
192 velop synthetic geochemical logs for pre-salt reservoirs in Brazil (de Oliveira
193 and de Carvalho Carneiro, 2021). Seven petrophysical logs: natural gamma-
194 ray, gamma-ray spectroscopy, density, photoelectric factor, neutron poros-
195 ity, nuclear magnetic resonance, and sonic formed the input variables.
196 The chemical element concentrations for Al, Ca, Fe, Mg, Na, Si, S, and
197 Ti were the prediction objectives. In addition to showing the best results,
198 AdaBoost was found to be the most practical tree-ensemble algorithm to
199 apply as it involved simpler pre-processing and control variable optimiza-
200 tion.

201 **3. Application Focus**

202 The objective of the ensemble models developed is to predict subsur-
203 face fluid flow characteristics. This is of interest to both engineers and sci-

204 entists with respect to, for example, the flow of fluids into producing wells
 205 penetrating oil and gas reservoirs, the seepage of wastewater through soil,
 206 and land subsidence as a consequence of groundwater and oil and gas ex-
 207 traction. With respect to oil and gas reservoirs, the main goal is to predict
 208 the performance of reservoirs at any future point in time and to optimize
 209 the petroleum fluid recovery under different operating conditions. Fluid
 210 flow in petroleum reservoirs is typically modelled using a set of non-linear
 211 PDEs. In general, these equations can be solved analytically (exact solu-
 212 tions) or numerically (approximate solutions).

213 A mixed GMsFEM method, as a numerical method has recently been
 214 proposed to solve Darcy's flow conditions (linear pressure gradient versus
 215 velocity) considering single-phase fluids in a porous medium character-
 216 ized by heterogeneties in two dimensions (i.e., matrix composition and
 217 fracture distribution) (Chen et al., 2020). The model approximates reser-
 218 voir pressure in multiscale space. It does so by applying several multiscale
 219 basis functions to a single coarse grid of the reservoir volume. The fluid
 220 velocity is directly estimated across a fine grid space.

The fluid flow conditions are defined as:

$$k^{-1}u + \nabla p = 0 \quad \text{in } \Omega \quad (4)$$

$$\nabla \cdot u = f \quad \text{in } \Omega \quad (5)$$

Heterogeneous boundary conditions are included:

$$u \cdot n = g \quad \text{on} \quad \partial\Omega \quad (6)$$

221 in which k is permeability, u is the Darcy velocity, p is the pressure, f is
 222 the source term, g is the normal to the Darcy velocity prevailing at the
 223 reservoir boundary, Ω is the computational domain and n is the outward
 224 unit norm vector on the boundary.

To illustrate the general solution framework of the mixed GMsFEM, τ^H is considered a confirming partition of Ω into finite elements with a coarse block size H , and τ^h is the fine grid partition with mesh size h . Assuming $V = H(\text{div}, \Omega)$ and $W = L^2(\Omega)$, the mixed finite element spaces become:

$$V_h = \left\{ v_h \in V : v_h(t) = (b_t x_1 + a_t, d_t x_2 + c_t), a_t, b_t, c_t, d_t \in \mathbb{R}, t \in \tau^h \right\}$$

$$W_h = \left\{ w_h \in W : w_h \text{ is a constant on each element in } \tau^h \right\}$$

$\{\Psi_j\}$ represents a set of multiscale base functions related to the coarse element. The multiscale space relating to pressure (p) can then be expressed as the linear extent of the local basis functions. This relationship is expressed as:

$$W_H = \oplus \{\Psi_i\} \quad \text{in} \quad \tau^H$$

In that form, the mixed GMsFEM is configured to find $(u_H, p_H) \in (V_h, W_H)$

constrained by:

$$\int k^{-1} u_H \cdot v_H - \int \operatorname{div}(v_H) p_H = 0 \quad \forall v_H \in V_h^0 \quad (7)$$

$$\int \operatorname{div}(u_H) w_H = \int f w_H \quad \forall w_H \in W_H \quad (8)$$

225 in which $u_H \cdot n = g_H$ on $\partial\Omega$ is relating to the coarse edges at the boundaries,
 226 whereas g_H is the average of function g at those coarse edges.

It is necessary to establish a multiscale space, W_H , to approximate p . This is achieved by solving local cell conditions for each coarse grid element by applying Dirichlet's boundary conditions. If $T_i \in \tau^H$ represents the coarse grid elements relating to Ω , the purpose is to find $(u_j^{(i)}, p_j^{(i)}) \in (V_h, W_h)|_{T_i}$ by solving the following problem on T_i :

$$k^{-1} u_j^{(i)} + \nabla p_j^{(i)} = 0 \quad \text{in } T_i \quad (9)$$

$$\operatorname{div}(u_j^{(i)}) = 0 \quad \text{in } T_i \quad (10)$$

227 where $(V_h, W_h)|_{T_i}$ is the restriction of (V_h, W_h) on T_i .

The coarse grid boundary element represents the junction of fine grid edges, i.e., $\partial T_i = \bigcup_{j=1}^{J_i} e_j$ in which J_i is the total number of fine grid edges at boundary T_i . $\delta_j^{(i)}$ represents a piecewise constant related to ∂T_i and the fine grid and $= 1$ for e_j and $= 0$ for the remaining fine grid edges.

Therefore, the boundary condition on the boundary of T_i is taken as the Dirichlet boundary condition:

$$p_j^{(i)} = \delta_j^{(i)} \quad \text{on } \partial T_i \quad (11)$$

By combining the local problem solutions a snapshot of spatial conditions is derived. Assuming $\Psi_j^{i,snap} := p_j^{(i)}$ defines the snapshot fields, then the snapshot space can be expressed:

$$W_{snap} = span \left\{ \Psi_j^{i,snap} : 1 \leq j \leq J_i, 1 \leq i \leq N_t \right\} \quad (12)$$

In the case of using the single-index notation:

$$W_{snap} = span \left\{ \Psi_i^{snap} : 1 \leq i \leq M_{snap} \right\} \quad (13)$$

228 where $M_{snap} = \sum_{i=1}^{N_t} J_i$ represents the total number of snapshot fields.

The snapshot space can then be further reduced by solving local grid problems. The local problem solutions are referred to as the offline space. The snapshot space corresponding to T_i becomes:

$$W_{snap}^{(i)} = span \left\{ \Psi_j^{i,snap} : 1 \leq j \leq J_i \right\}$$

In a local grid problem, the real number $\lambda \geq 0$ and the function $p \in W_{snap}^{(i)}$

need to be derived

$$a_i(p, w) = s_i(p, w) \quad \forall w \in W_{snap}^{(i)} \quad (14)$$

For each T_i :

$$a_i(p, w) = \sum_e k[p][w] \quad \text{and} \quad s_i(p, w) = \int kpw \quad \text{in } T_i \quad (15)$$

229 in which $[p]$ and $[w]$ are the jump of functions p and w , respectively. Also,
230 e represents the fine edge interior of T_i .

The eigenvalues of Equation 14 are arranged in increasing order:

$$\lambda_1^{(i)} < \lambda_2^{(i)} < \dots < \lambda_{l_i}^{(i)} \quad (16)$$

where $\lambda_k^{(i)}$ denotes the k th eigenvalue for T_i . The corresponding eigenvectors are $Z_k^{(i)} = (Z_{kj}^{(i)})_{j=1}^{J_i}$ with $Z_{kj}^{(i)}$ being the j th component of the vector $Z_k^{(i)}$. Initial l_i eigenfunctions are selected to represent the offline space. Offline basis functions are then defined as:

$$\Psi_k^{i,off} = \sum_{j=1}^{J_i} Z_{kj}^{(i)} \Psi_j^{i,snap} \quad k = 1, 2, \dots, l_i.$$

Then, global offline space becomes:

$$W_{off} = \text{span} \left\{ \Psi_k^{i,off} : 1 \leq k \leq l_i, 1 \leq i \leq N_t \right\}$$

Applying single-index notation, the global offline space can be defined as:

$$W_{off} = span \left\{ \Psi_k^{off} : 1 \leq k \leq M_{off} \right\}$$

231 where $M_{off} = \sum_{i=1}^{N_t} l_i$ is the total number of offline basis functions.

Each Ψ_k^{off} can be expressed by a vector ψ_k^{off} which contains coefficients from Ψ_k^{off} relating to the fine grid basis functions. Thus:

$$R_{off} = \left[\psi_1^{off}, \dots, \psi_{M_{off}}^{off} \right]$$

The offline space is mapped using these functions to the fine grid space. The mixed GMsFEM system (Equations 7 and 8) is expressed in matrix terms as:

$$M_{fine} U_H + B_{fine}^T R_{off} P_H = 0 \quad (17)$$

$$R_{off}^T B_{fine} U_H = R_{off}^T F_H \quad (18)$$

232 M_{fine} constitutes a symmetric, positive definite, and sparse matrix. U_H
 233 and P_H are the unknown fluid velocity and pressure vectors that describe
 234 grid spaces V_h and W_H , respectively. Execution of the mixed GMsFEM
 235 therefore requires two fine grid matrices to be constructed (M_{fine} , B_{fine})
 236 accompanied by one offline matrix (R_{off}).

Fluid velocity can be solved directly from the fine grid matrix combination. Considering k as a diagonal tensor, M_{fine} is readily estimated with

the diagonal matrix \widehat{M}_{fine} , applying the trapezoidal quadrature rule. The convergence rate of that easier-to-execute system is essentially the same as that using the unmodified matrix. M_{fine} can therefore be replaced by that diagonal matrix \widehat{M}_{fine} without compromising prediction accuracy. As \widehat{M}_{fine} is easier to invert, the system described by Equations 17 and 18 is solved as follows:

$$-R_{off}^T B_{fine} \widehat{M}_{fine}^{-1} B_{fine}^T R_{off} P_H = R_{off}^T F_H$$

237 Taking this approach, an original mixed formulation is expressed approx-
 238 imately by a positive-definite, sparse linear system. In that linear system,
 239 fewer pressure unknowns are involved for each coarse-grid element.

240 Generally, the number of PDEs requiring solutions to enable multiscale
 241 basis functions to be derived is dependent on the number of local cell and
 242 local eigenvalue problems involved. The local cell problem relating to
 243 the coarse grid relates to the original system definition but excludes the
 244 source function in Equation 5. A boundary condition (δ) relates to the
 245 coarse grid boundary; $\delta=1$ for fine grid edges and $\delta=0$ for coarse
 246 grid edges. Local cell problems are therefore determined by the fine grid
 247 edges impacting the coarse grid boundary. In the model configured for
 248 this study, the number of fine grid edges/coarse grid boundary is 12.

249 In this study, the Karhunen-Loeve expansion was used to parameterize
 250 the heterogeneous permeability field. This Gaussian random field genera-

251 tion method decomposes a random process into the eigenvalue and eigen-
252 function of its covariance kernel. The fine grid system adopted involves a
253 uniform 30×30 mesh. On the other hand, a sparser, uniform 10×10 mesh
254 was applied to represent the coarse grid network (Figure 1). This configu-
255 ration consists of 1300 separate PDEs, made up of 1200 PDEs addressing
256 the local cell problems (100 coarse grid mesh units by 12 fine grid edges)
257 plus 100 local eigenvalue problems (one per each 100 coarse grid mesh
258 units). The input to this model is comprised of a randomly-generated per-
259 meability field. For each permeability field, there are five basis functions,
260 (numbered Basis 1 to 5). Basis 1 is a piecewise constant, with binary val-
261 ues of -1 and $+1$. Basis 1 is defined as part of the finite element method,
262 it therefore requires no training for DL modeling. On the other hand, Ba-
263 sis 2 to 5 take values distributed across the range $(-1, +1)$, and therefore
264 require training for DL modeling.

265 **4. Stacking CNN Ensemble Model**

266 A schematic of the proposed stacking CNN ensemble model is given
267 in Figure 2. The data generation and pre-processing steps involved are
268 described in Sub-section 4.1. The mechanism adopted for training the
269 base learners is outlined in Sub-section 4.2. The procedure for combining
270 the DL model results is presented in Sub-section 4.3.

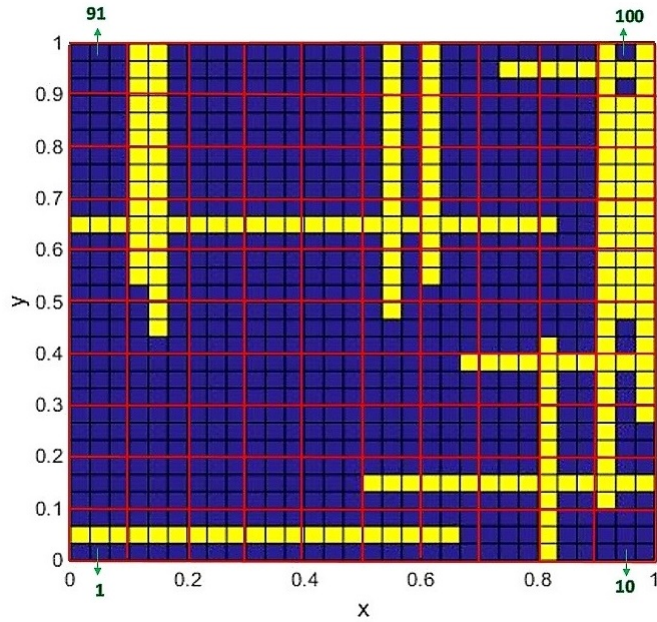


Figure 1: Schematic description of the permeability field of a simulated fractured porous reservoir formation. Matrix permeability is assumed to be 4 milliDarcies (mD). Fracture permeability is assumed to be 2000 mD. Fine grid squares represent the formation matrix (blue) in some cases and fractures (yellow) in other cases (selected randomly). The red lines define the coarse grid. Each coarse grid square contains of nine fine grid squares. There are fifteen fractures assigned to this porous medium.

271 *4.1. Data preparation and pre-processing*

272 The ranges of permeability values applied to the formation matrix were
 273 1, 2, 3, 4, and 5 mD, and to the fractures were 500, 750, 1000, 1250, 1500,
 274 1750, and 2000 mD. The number of fractures per 10×10 coarse grid was
 275 varied between 1 and 25. This ranges meant that 875 cases (5 matrix per-
 276 meabilities by 7 fracture permeabilities by 25 different fracture densities)
 277 in total required evaluation by the DL model. The format defined for the
 278 permeability field was as a vector (900×1), subsequently adjusted to be

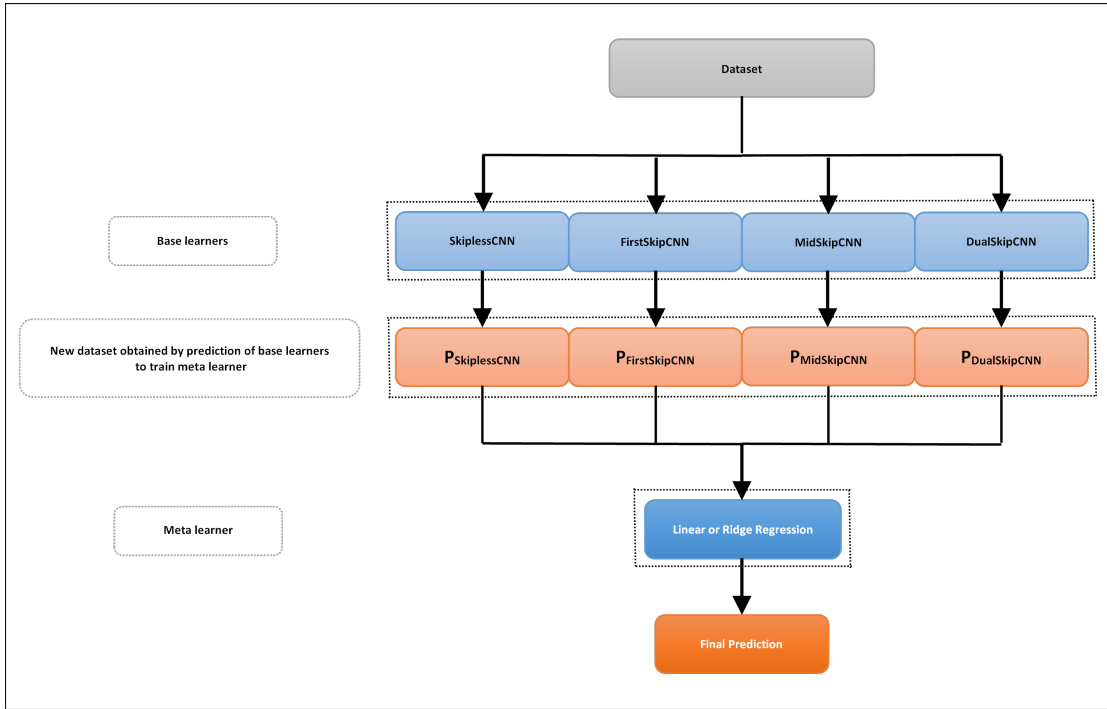


Figure 2: The workflow diagram of the stacking CNN ensemble model. Four base learners of SkiplessCNN, FirstSkipCNN, MidSkipCNN, and DualSkipCNN are developed using training/validation subsets. After being trained, they are used to make predictions on the validation data ($P_{\text{SkiplessCNN}}$, $P_{\text{FirstSkipCNN}}$, $P_{\text{MidSkipCNN}}$, and $P_{\text{DualSkipCNN}}$). Then, two meta models are separately developed using linear regression and ridge regression. Once the meta models are trained, they can be used to make predictions on the testing data (28,879 samples).

279 expressed as a 2D tensor (100×9), in which, coarse grid units=100 and
 280 each coarse grid contains 9 fine grids. Each row in the array therefore rep-
 281 represents a coarse grid. Such a configuration enables the use of 2D CNN
 282 kernels. However, it was necessary to maintain the five basis functions as
 283 900×1 vectors, so that they could be evaluated in the FC layers forming
 284 the final section of the CNN network.

285 Each of the 875 cases was evaluated 350 times as part of model train-

286 ing. Additionally, 40 validation runs, and 40 independent testing runs
287 were executed for each case. Matlab code was used to run the mixed GMS-
288 FEM models. These models generated 376,250 data records one for each
289 10×10 coarse grid configuration. 306,250 records were used for DL model
290 training, 35,000 records for DL model validation, and 35,000 for indepen-
291 dently testing the trained and validated models. The random generation
292 of each permeability field (10×10 coarse grid) involves the possibility that
293 some duplicate fields could be generated. Consequently, the generated
294 dataset was filtered to remove any duplicate data records. This is neces-
295 sary to remove the risk of introducing bias towards specific model configu-
296 rations in the DL analysis. This data filtering step removed 1739 duplicate
297 training data records , 579 duplicate validation data records, and 6121
298 duplicate testing data records were excluded. This pre-processing step re-
299 duced the training subset to 304,511 data records, the validation subset
300 to 34,421 data records, and the independent testing subset to 28,879 data
301 records.

302 4.2. *Training the base learners*

303 The CNN algorithm has, over recent years, become one of the most
304 trusted DL models with respect to many application domains (Rao et al.,
305 2017; Chen and He, 2018; Pratt et al., 2019; Hakim et al., 2022). The CNN
306 was originally designed to solve problems with 2D arrays, particularly
307 images, although it can also be applied to 1D arrays. It progressively and

308 flexibly learns feature relationships, spatially in 2D models, by applying
309 an optimizer of choice to various types of network layers. The primary
310 CNN layer types are (i) convolutional, (ii) pooling, and (iii) FC, usually
311 configured in that sequence. In mathematical sciences, convolution is a
312 specialized linear operation on two functions that gives a third modified
313 function. In the context of CNN, the fundamental idea is to consider an
314 input (an array of numbers) as the first function and a convolutional filter
315 (kernel) as the second. A kernel is a relatively small array of randomly
316 generated numbers. The kernel moves over the whole input. The dot
317 product of the kernel and input is calculated at each sub-region (with the
318 same size as the kernel) of the input, obtaining an output value in the
319 corresponding location of the convolved input. This process produces a
320 feature map and is performed using different kernels. The outputs of the
321 convolution process are passed through an activation (transfer) function.
322 Such functions typically transform a linear operation into a nonlinear sys-
323 tem (Yamashita et al., 2018; Elgendy, 2020).

324 The key difference between a parameter and a hyperparameter is that
325 a model's parameters are automatically updated during the training pro-
326 cess, whereas hyperparameters are set manually before the model begins
327 training, for example, the size and number of kernels. Including more
328 convolutional layers in a CNN model increases the number of parameters.
329 The more parameters there are in a model, the more computationally ex-
330 pensive the learning process is. This is where a subsampling operation can

331 be useful. In DL, pooling layers use statistical functions (maximum and
332 average pooling) to decrease the number of trainable parameters. This can
333 decrease the computational complexity of mathematical operations and
334 sometimes improve the robustness of feature maps. Pooling layers come
335 after convolutional layers (Yamashita et al., 2018; Elgendy, 2020).

336 When the output of the network is in the format of a vector, feature
337 maps in the final convolutional or pooling layer are first flattened to a one-
338 dimensional array, and then connected to FC layers. In FC layers (dense
339 layers), each neuron of a layer is connected to whole neurons in the previ-
340 ous layer and the next layer. It is common to put a dropout layer after each
341 FC layer (except the output layer) at the end of a CNN model. Dropout
342 omits a percentage of neurons in the previous FC layer. This percentage,
343 as a hyperparameter, is defined when constructing a network. During the
344 training process, some neurons may dominate, producing errors. Dropout
345 balances a network, checking that all neurons work equally to minimize
346 the cost function as much as possible (Yamashita et al., 2018; Elgendy,
347 2020).

348 To develop the base learner for this study, the 304,511 training data
349 records, together with the 34,421 validation data records were employed.
350 Distinct CNN model configurations, involving various combinations of
351 convolutional, pooling, FC, Batch Normalization (BN), regularization, and
352 dropout filtering were tested separately for each basis function requiring
353 training (Basis 2 to 5). A similar optimal CNN configuration for each of

354 those four basis functions (Figure 3: identified as SkiplessCNN) was ob-
 355 tained. That initial CNN architecture consists of five convolutional layers,
 356 two FC layers but does not include any pooling layers. Convolutional lay-
 357 ers 1 to 5 (CONV1 to CONV5) consist of 5, 10, 15, 20, and 25 kernels,
 358 respectively. To determine the size of a convolution output for an input
 359 with the size of $I_h(\text{height}) \times I_w(\text{width})$ and a kernel with the size of $K_h \times K_w$,
 360 we can use Equation 19 if the padding is set to ‘valid’:

$$\begin{aligned} \text{Output height} = O_h &= (I_h - K_h) / S_h + 1 \\ \text{Output width} = O_w &= (I_w - K_w) / S_w + 1 \end{aligned} \tag{19}$$

361 where S_h and S_w are the vertical and horizontal strides. When padding is
 362 set to ‘same’, the size does not change. The kernel size for all convolutional
 363 layers is 3×3 , and $S_h = S_w = 1$. The padding was set to ‘valid’ only for
 364 CONV5. This means there was no padding for the first four convolutional
 365 layers. Therefore, CONV1, CONV2, CONV3, CONV4, and CONV5 have
 366 the size of 98×7 , 96×5 , 94×3 , 92×1 , 92×1 , respectively.

367 Each convolutional layer is followed by a single BN layer of the same
 368 dimensions. Typically, neural network models converge more quickly when
 369 the input to each layer is normalized; hence the value of adding the BN lay-
 370 ers. Each FC layer contains 2000 neurons. For a given neuron or kernel,
 371 the inputs are multiplied by weights and the resulting products summed
 372 together. A bias term is then applied to that sum. Such rigid computations
 373 mean that only linear transformations are performed on the layer inputs

374 using the weights and biases to generate the layer outputs. Although this
375 operation makes the neural network simpler, it is less powerful and un-
376 able to learn complex patterns in a dataset. This is where the activation
377 function is beneficial. Mathematically, this can be represented as shown
378 in Equation 20 where w_i represents the weight value, z_i is the input value,
379 b is the bias, f refers to the activation function applied, and y is the de-
380 pendent variable prediction output. The developed models in this study
381 used the ‘Rectified Linear Unit (ReLU)’ activation function for the convo-
382 lutional layers, ‘sigmoid’ activation function for the FC layers, and ‘linear’
383 activation function for the output.

$$y = f\left(\sum_{i=1}^n (w_i z_i) + b\right) \quad (20)$$

384 In order to better understand the standard architecture (i.e., Skipless-
385 CNN) developed in this study, it is compared to structurally similar CNN
386 architectures AlexNet (Krizhevsky et al., 2017) and VGGNet, also known
387 as VGG16 (Simonyan and Zisserman, 2014). AlexNet has five convolu-
388 tional layers, three of which are followed by maximum pooling layers to
389 decrease the computational cost. The number of kernels in each con-
390 volutional layer is 96, 256, 384, 384, and 256. There are two FC lay-
391 ers of 4096 neurons and a 1000-neuron output layer at the end of the
392 network. VGGNet contains thirteen convolutional layers, five maximum
393 pooling layers, two FC layers of 4096 neurons, and an output layer with

394 1000 neurons. The number of kernels used in the convolutional sections
395 is 64, 128, 256, and 512. Similar to common CNN architectures, going
396 deeper through the structure of developed models, the number of feature
397 maps increases and their size decreases. However, the number of feature
398 maps (equals the filters number) defined in this research is significantly
399 less than that of common CNN models. In DL, pooling layers are primar-
400 ily used to decrease the number of trainable parameters, mostly when the
401 input shape is high, e.g., in AlexNet whose input shape is 224×224 . How-
402 ever, the input dimension in this research is 100×9 . This is why there
403 is no pooling layer in our developed models. BN, similar to AlexNet, has
404 helped to prevent over-fitting. As with AlexNet and VGGNet, the num-
405 ber of neurons (units) remained constant in FC layers, but no drop out
406 layer was used in the proposed structure because it had a negative effect
407 on the performance. The base structure of this work is for a regression-
408 type problem, while AlexNet and VGGNet were essentially designed for a
409 classification intent. Therefore, a linear activation function is used in our
410 model for the output layer, but a softmax in AlexNet and VGGNet.

411 The CNN training process seeks to find optimum values for weights
412 and biases applied to kernels (convolutional layers) and neurons (FC lay-
413 ers). Such values generate the lowest collective errors for all data records
414 evaluated between actual and predicted dependent variable values. The
415 back-propagation algorithms are commonly applied to train many types
416 of neural network. They calculate the gradient of the loss function (cost

417 function) using the values assigned to the weights and biases. The loss
418 function is a measure of how well an algorithm models a training dataset
419 by evaluating the similarity between real and predicted outputs.

420 The different optimizers available all strive to achieve a minimum loss
421 or cost value. Multi-layer perceptron neural networks focus on the feed-
422 forward sequence through its layered structure on to which weights and
423 biases are initialized. However, in training the backward pathway is used
424 to modify the layer weights and biases in each iteration. In that way back
425 propagation acts to improve a model's performance.

426 The CNN models in this study were constructed using Keras with Ten-
427 sorFlow as a backend on Python. The models were compiled using 'MSE'
428 as the loss (objective) function. The learning rate is a key DL hyper-
429 parameter. It states how quickly a model learns in each epoch that pa-
430 rameters are updated. When it is too small, the training process takes a
431 long time. If too large, it results in sub-optimal CNN learning, locking
432 into sets of weights and biases too quickly, which can lead to a less stable
433 training process that tends to converge prematurely. Hence, setting the
434 right value of the learning rate is crucial. Adaptive methods such as Adam
435 can be used to automatically resolve this issue. Adam applies distinctive
436 learning rates to each scalar variable. It progressively adapts those rates
437 throughout the training iterations, with those adaptations being influenced
438 by partial-derivative trends of rates applied to each variable in previous
439 model iterations. Adam is gradient based in its calculations and bene-

440 fits from a combination of its AdaGrad component to cope with sparse
441 gradients and an RMSProp function in its application. It is suitable for
442 DL applications to large datasets with many data records and/or multiple
443 variables. The Adam learning rate can adjust at a finer scale as the opti-
444 mum values are approached, although in some cases such fine tuning can
445 result in overfitting. AMSGrad extends the performance of the Adam op-
446 timizer by converging in a more effective and smoother manner, avoiding
447 step changes. By storing the highest values of second- momentum vectors
448 generated in all previous model iterations, AMSGrad is able to normal-
449 ize the moving average gradient in each iteration. To benefit from these
450 advantages, the AMSGrad with the default values i.e., the initial global
451 learning rate = 0.001, $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 1e - 7$ has
452 been applied to the CNN models developed for this study. In addition, the
453 models were trained with a batch size of 32 samples over 100 epochs.

454 In feed-forward neural networks with multiple layers, such as most DL
455 models, back propagation works from the latter layers back through mul-
456 tiple layers to reach the initial layers. This extended sequence can result
457 in the gradient being reduced rapidly, in very few model iterations, to a
458 value close to zero. This generally unfavorable premature convergence is
459 referred to the “vanishing gradient phenomenon”. It is ineffective because
460 it prematurely halts the training process before the early layers of the net-
461 work have fully explored potentially more favorable values. A beneficial
462 strategy that acts to reduce the risk of premature convergence is to involve

463 “skip” connections between certain network layers, acting as short circuits
464 for the back-propagation sequence (He et al., 2016). The introduction of
465 skip connections enables the gradient to be directly back propagated to
466 earlier layers of a CNN. Skip connections (shortcuts) are involved in three
467 of the base learners used in this study (Figure 3). How and where in the
468 CNN structure the shortcuts are located differs from scheme to scheme:

- 469 1. **FirstSkip**: a single skip connection from the first convolutional layer
470 to the last one.
- 471 2. **MidSkip**: a single skip connection from the middle convolutional
472 layer to the last layer.
- 473 3. **DualSkip**: two skip connections from the middle convolutional layer
474 to the last and the second-to-last layers.

475 **FirstSkip** adds a single shortcut from the output of the first convo-
476 lution layer to the last convolutional block. The input and output of this
477 part have the same dimension of 98×7 because an identity type of shortcut
478 is used. **MidSkip** is designed to discover how much a shortcut from the
479 middle layer to the final layer can improve the performance of a model.
480 Here, the input and output of this section with the shortcut have a dimen-
481 sion of 94×3 . **DualSkip** was developed mainly to gain knowledge of the
482 effect of involving the raw input features along with **MidSkip**. In all three
483 cases, the main path and the shortcut meet each other before applying
484 the activation function. For all three architectures, the FC layers remain
485 unchanged.

486 Adding the skip connections increases the complexity of the base struc-
487 ture, in terms of the number of parameters. There are as many as 10,414,170
488 trainable and 150 non-trainable parameters for the base structure, with-
489 out skip connections. By adding FirstSkip, the number of trainable and
490 non-trainable parameters increases to 12,670,510 and decreases to 110,
491 respectively. For MidSkip, there are 14,272,340 trainable and 130 non-
492 trainable parameters. For DualSkip, the number of trainable and non-
493 trainable parameters changes to 14,270,975 and 120, respectively.

494 4.3. *Combination of the base learner outputs*

495 A new training dataset for the meta learner was established by pro-
496 viding all data records from the validation subsets to each of the four
497 sub-models and collecting the predictions they generated i.e., $P_{\text{SkiplessCNN}}$,
498 $P_{\text{FirstSkipCNN}}$, $P_{\text{MidSkipCNN}}$, and $P_{\text{DualSkipCNN}}$. These resulted in four (refer-
499 ring to the number of base learners) arrays with the shape [34421, 900],
500 the first element referring to the number of validation data and the second
501 to the output (basis function). Thus, a 3D array was developed with the
502 shape [34421, 4, 900], which was transformed into a [34421, 3600] shaped
503 array. This flattened input data, along with their output was used to train
504 a meta learner.

505 As mentioned earlier, there are 875 cases that each need to be pro-
506 cessed through the sequence of training, validation, and testing. Given
507 that the input/output dimensions are so large, it did not make sense to ap-

508 ply boosting to focus on samples in the dataset that have been predicted
509 incorrectly by the previous models in the sequence. Bagging is usually
510 applied to relatively small datasets. Additionally, conducting bootstrap
511 sampling incorporating all 875 cases was not feasible for addressing this
512 large dataset. Therefore, stacking was selected to establish the ensemble
513 model. A stacked generalization method was chosen, mainly because it is
514 more flexible mathematically than voting or weighted average methods.
515 To be more specific, the four base learners were combined into an ensem-
516 ble model using linear and ridge regression, separately.

517 **5. Evaluation**

518 The prediction errors associated with each CNN model developed are
519 assessed using two statistical error metrics: the R^2 and MSE. R^2 values
520 can exist within a range $-\infty$ and 1, with values closest to 1, representing
521 the better prediction performance. MSE, by definition, has to be a non-
522 negative value, and where values closer to zero represent the better perfor-
523 mance. Table 1 presents the prediction error results for the base learners
524 using the training and validation subsets. All the constructed Skipless-
525 CNN models yield satisfactory results for the training samples; the best is
526 for Basis 4 with an R^2 of 0.9156 and MSE of 0.0126. All three skip connec-
527 tion schemes enhance performance over the standard structure for all ba-
528 sis functions evaluated (i.e., Basis 2, 3, 4, and 5) by their training subsets.
529 The defined schemes have the maximum effect on the base model for Ba-

530 sis 5 and the minimum effect for Basis 4. For example, R^2 increases from
 531 0.8466 to 0.8844, 0.9026, and 0.8847 for Basis 5 by including FirstSkip,
 532 MidSkip, and DualSkip into the standard structure, respectively. MidSkip
 533 and DualSkip perform marginally better than FirstSkip.

Table 1: Prediction error analysis of the base learners applied to training and validation data subsets.

Subset	Model	R^2				MSE			
		Basis 2	Basis 3	Basis 4	Basis 5	Basis 2	Basis 3	Basis 4	Basis 5
Training	SkiplessCNN	0.8657	0.8952	0.9156	0.8466	0.0327	0.0220	0.0126	0.0100
	FirstSkipCNN	0.8908	0.9219	0.9247	0.8844	0.0266	0.0164	0.0112	0.0075
	MidSkipCNN	0.9083	0.9302	0.9372	0.9026	0.0224	0.0147	0.0093	0.0063
	DualSkipCNN	0.9002	0.9327	0.9283	0.8847	0.0243	0.0141	0.0107	0.0075
Validation	SkiplessCNN	0.7770	0.8237	0.8777	0.7816	0.0544	0.0371	0.0182	0.0142
	FirstSkipCNN	0.7814	0.8160	0.8798	0.7974	0.0529	0.0387	0.0181	0.0132
	MidSkipCNN	0.7867	0.8139	0.8802	0.8160	0.0519	0.0391	0.0179	0.0120
	DualSkipCNN	0.7900	0.8434	0.8811	0.8038	0.0512	0.0329	0.0176	0.0128

534 The prediction error performance of the SkiplessCNN models is ac-
 535 ceptable for the validation data subsets, with an R^2 of 0.7770 to 0.8777,
 536 and MSE of 0.0142 to 0.0544. FirstSkip has a marginally positive effect on
 537 the validation subset with respect to Basis 2, 4, and 5 models. For instance,
 538 regarding Basis 5, the R^2 value increases from 0.7816 to 0.7974 and MSE
 539 decreases from 0.0142 to 0.0132. However, it has an adverse effect on the
 540 Basis 3 model. Specifically, R^2 decreases from 0.8237 to 0.8160. Com-
 541 pared to FirstSkip, MidSkip has a more positive effect on the Basis 2, 4,

542 and 5 models. Furthermore, it has a negative effect on the Basis 3 model.
543 DualSkip is beneficial in all cases related to validation samples, especially
544 for Basis 3 and 5. For example, for Basis 5, R^2 increases from 0.7816 to
545 0.8038.

546 The results obtained for FirstSkip imply that transferring feature maps
547 from earlier convolutional layers to final ones has a very positive effect on
548 the training dataset. This architecture has a marginally positive impact on
549 the validation subset for Basis 2, 4, and 5 models, but an adverse impact
550 on the Basis 3 model. In other words, the corresponding skip connection
551 tends to make the predictive model focus more on capturing the underlying
552 trend of the training (seen) subset.

553 Compared to FirstSkip, flowing information from the middle convolu-
554 tional layer to the last layer via the MidSkip skip connection has a more
555 positive impact on all basis functions models of the training subset and
556 the Basis 2, 4, and 5 models of the validation subset. This suggests that
557 the feature maps of the middle convolution process contain important in-
558 formation.

559 Adding two simultaneous skip connections (DualSkip) favorably af-
560 fects all basis functions with respect to the training and validation sub-
561 sets. By comparing the architectures and results produced using MidSkip
562 and DualSkip, the positive role of transferring raw feature maps is under-
563 standable. Therefore, enriching the last convolutional blocks with infor-
564 mation hidden in the neighboring layers is more efficient than enriching

565 them using earlier convolutional blocks near the input layer.

566 Figures 4 and 5 illustrate whether the combined models significantly
567 influence the performance of the base learners based on the testing subset.
568 The trend for the base learners over the testing subset is the same as their
569 performance in the validation subset. The most likely reason is that the
570 validation and testing subsets were drawn from the same data distribu-
571 tion. Both subsets also consisted of the same number of data records. It
572 is apparent from Figures 4 and 5 that the ensemble models built by either
573 linear or ridge regression perform substantially better on the testing sub-
574 set than the individual models. In the case of applying linear regression,
575 the R^2 and MSE lie in the range of 0.8456 to 0.9191, and 0.0092 to 0.0369,
576 respectively. The results reveal that ridge regression works marginally bet-
577 ter than linear regression with an R^2 ranging from 0.8539 to 0.9220, and
578 ranging from MSE of 0.0090 to 0.0349.

579 Figure 6 presents an example of the pressure distributions for a repre-
580 sentative permeability field, whose matrix permeability and fracture per-
581 meability are 1 and 1750 mD, respectively. The figure displays a very close
582 match between the actual pressure distribution and the one obtained by
583 applying the predictions derived from the developed ridge regression en-
584 semble model in this study. The reservoir pressure distribution is impor-
585 tant information used to determine the potential recovery factor of oil/gas
586 reservoirs. These results confirm that this study's innovative approach to
587 develop efficient and reliable DEL models to assist this specific aspect of

588 fluid-flow simulation is successful and worthy of further development.

589 **6. Discussion**

590 The linear regression is one of the most straightforward approaches
591 to predict output via a linear function of input features. In the context
592 of ML, it refers to the most usual least square linear regression method
593 that attempts to minimize the cost function. A drawback, however, is that
594 it does not penalize high magnitude weights in its error function and it
595 assumes independence between its features. These characteristics can lead
596 in some cases to certain features being assigned very high weights during
597 the training. The cost function for linear regression is typically expressed
598 as:

$$cost\ function_{linear} = \sum_{i=1}^m (\mathbf{Y} - \hat{f}(\mathbf{X}))^2 \quad (21)$$

599 The ridge regression, as a modification of linear regression, involves
600 a penalty (L2 regularization) to its error term, calculated as the sum of
601 squared value of the weights. Giving a penalty in such a way results in a
602 set of more evenly distributed weights. The cost function for ridge regres-
603 sion becomes:

$$cost\ function_{ridge} = \sum_{i=1}^m (\mathbf{Y} - \hat{f}(\mathbf{X}))^2 + \alpha \sum_{j=1}^p (w_j)^2 \quad (22)$$

604 Here, α is included as a coefficient to penalize weights. It can take

605 different values. A ridge model with $\alpha = 0$ is the same as a simple linear
606 regression. As the α value nears infinity, an increasing number of coef-
607 ficients of the model becomes zero until it is just a flat model with an
608 intercept. In this study, we used the default value of "one" for all cases.

609 It was expected that ridge regression would perform better than linear
610 regression, and the results presented in Section 5 confirm this point.

611 Analysis of the case study, the proposed method, and the results lead to
612 three recommendations for future research. First, the scope of the present
613 study was restricted to 2D porous media with vertical and horizontal frac-
614 tures. It is recommended to extend it to 3D porous media and incorporate
615 inclined fractures. Moreover, it would be worthwhile to consider a wider
616 range of permeability values for both the matrix and fractures. These ad-
617 justments would provide a more comprehensive representation of subsur-
618 face conditions. Second, the effectiveness of the stacked generalization
619 method using linear regression and ridge regression was confirmed on the
620 given task. However, there is still room for improvement in the testing
621 subset. To further enhance performance, exploring more advanced en-
622 semble techniques could be one direction to consider. Third, developing
623 more diverse base learners could also provide valuable insights for further
624 improving the performance of the stacking ensemble.

625 **7. Conclusions**

626 The substantial quantity of parameters involved in DL means that a
627 large number of samples must be processed to provide effective results.
628 The ReLU has emerged as a popular activation functions applied in DL,
629 and AMSGrad, an enhanced version of the Adam optimizer, improves
630 DL convergence. By using skip connection (shortcut) schemes during
631 gradient-based training, such as back propagation, the vanishing gradient
632 problem can be mitigated. These four DL performance features were ap-
633 plied to a case study to predict subsurface fluid flow and simplify a time-
634 consuming component of oil/gas reservoir simulation. For this purpose,
635 four distinct CNN learners - SkiplessCNN, FirstSkipCNN, MidSkipCNN,
636 and DualSkipCNN - were developed for each multiscale basis function.
637 Linear regression and ridge regression were then used separately to com-
638 bine the four CNN into an ensemble model. The results confirm the effec-
639 tiveness of the two tested ensemble architectures since they strike a more
640 stabilized balance between bias and variance errors.

641 **Acknowledgment**

642 We appreciate the University of Liverpool for providing us with a pow-
643 erful computer system to run our codes. Also, we would like to thank Ms.
644 Roslyn Joy Irving and Dr. Trevor Mahy for providing helpful comments
645 and suggestions.

646 **Funding**

647 This work is partially supported by Key Program Special Fund in XJTLU
648 (KSF-E-50), XJTLU Research Development Funding (RDF-19-01-15), and
649 Research Grants of School of Mathematics and Physics.

650 **References**

651 Breiman, L., 1996a. Bagging predictors. *Machine learning* 24, 123–140.

652 Breiman, L., 1996b. Stacked regressions. *Machine learning* 24, 49–64.

653 Chen, J., Chung, E.T., He, Z., Sun, S., 2020. Generalized multiscale ap-
654 proximation of mixed finite elements with velocity elimination for sub-
655 surface flow. *Journal of Computational Physics* 404, 109133.

656 Chen, S., He, H., 2018. Stock prediction using convolutional neural net-
657 work, in: *IOP Conference series: materials science and engineering*, IOP
658 Publishing. p. 012026.

659 Elgendy, M., 2020. *Deep learning for vision systems*. Simon and Schuster.

660 Hakim, W.L., Rezaie, F., Nur, A.S., Panahi, M., Khosravi, K., Lee, C.W.,
661 Lee, S., 2022. Convolutional neural network (cnn) with metaheuristic
662 optimization algorithms for landslide susceptibility mapping in icheon,
663 south korea. *Journal of environmental management* 305, 114367.

- 664 He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image
665 recognition, in: Proceedings of the IEEE conference on computer vision
666 and pattern recognition, pp. 770–778.
- 667 Ho, T.K., 1998. The random subspace method for constructing decision
668 forests. IEEE transactions on pattern analysis and machine intelligence
669 20, 832–844.
- 670 Kanda, E., Epureanu, B.I., Adachi, T., Tsuruta, Y., Kikuchi, K., Kashihara,
671 N., Abe, M., Masakane, I., Nitta, K., 2020. Application of explainable
672 ensemble artificial intelligence model to categorization of hemodialysis-
673 patient and treatment using nationwide-real-world data in japan. Plos
674 one 15, e0233491.
- 675 Kardani, N., Zhou, A., Nazem, M., Shen, S.L., 2021. Improved prediction
676 of slope stability using a hybrid stacking ensemble method based on
677 finite element analysis and field data. Journal of Rock Mechanics and
678 Geotechnical Engineering 13, 188–201.
- 679 Krizhevsky, A., Sutskever, I., Hinton, G.E., 2017. Imagenet classification
680 with deep convolutional neural networks. Communications of the ACM
681 60, 84–90.
- 682 Kumari, P., Toshniwal, D., 2021. Extreme gradient boosting and deep neu-
683 ral network based ensemble learning approach to forecast hourly solar
684 irradiance. Journal of Cleaner Production 279, 123285.

- 685 Liu, Z., Gilbert, G., Cepeda, J.M., Lysdahl, A.O.K., Piciullo, L., Hefre, H.,
686 Lacasse, S., 2021. Modelling of shallow landslides with machine learn-
687 ing algorithms. *Geoscience Frontiers* 12, 385–393.
- 688 Nguyen, M., Logofătu, D., 2018. Applying tree ensemble to detect
689 anomalies in real-world water composition dataset, in: *International*
690 *Conference on Intelligent Data Engineering and Automated Learning*,
691 Springer. pp. 429–438.
- 692 Oh, H.J., Syifa, M., Lee, C.W., Lee, S., 2019. Land subsidence suscepti-
693 bility mapping using bayesian, functional, and meta-ensemble machine
694 learning models. *Applied Sciences* 9, 1248.
- 695 de Oliveira, L.A.B., de Carvalho Carneiro, C., 2021. Synthetic geochem-
696 ical well logs generation using ensemble machine learning techniques
697 for the brazilian pre-salt reservoirs. *Journal of Petroleum Science and*
698 *Engineering* 196, 108080.
- 699 Pratt, H., Williams, B., Broadbent, D., Harding, S., Coenen, F., Zheng, Y.,
700 2019. Learning the features of diabetic retinopathy with convolutional
701 neural networks. *EUROPEAN JOURNAL OF OPHTHALMOLOGY* 29,
702 NP15–NP16.
- 703 Rao, Y., He, L., Zhu, J., 2017. A residual convolutional neural network for
704 pan-shaprening, in: *2017 International Workshop on Remote Sensing*
705 *with Intelligent Processing (RSIP)*, IEEE. pp. 1–4.

- 706 Sahin, E.K., Demir, S., 2023. Greedy-automl: A novel greedy-based stack-
707 ing ensemble learning framework for assessing soil liquefaction poten-
708 tial. *Engineering Applications of Artificial Intelligence* 119, 105732.
- 709 Schapire, R.E., 1990. The strength of weak learnability. *Machine learning*
710 5, 197–227.
- 711 Simonyan, K., Zisserman, A., 2014. Very deep convolutional networks for
712 large-scale image recognition. *arXiv preprint arXiv:1409.1556* .
- 713 Tüysüzöğlü, G., Birant, D., 2020. Enhanced bagging (ebagging): A novel
714 approach for ensemble learning. *International Arab Journal of Informa-
715 tion Technology* 17.
- 716 Verma, R., Chandra, S., 2023. Repute: A soft voting ensemble learning
717 framework for reputation-based attack detection in fog-iot milieu. *En-
718 gineering Applications of Artificial Intelligence* 118, 105670.
- 719 Wang, G., Hao, J., Ma, J., Jiang, H., 2011. A comparative assessment of
720 ensemble learning for credit scoring. *Expert systems with applications*
721 38, 223–230.
- 722 Wang, G., Sun, J., Ma, J., Xu, K., Gu, J., 2014. Sentiment classification: The
723 contribution of ensemble learning. *Decision support systems* 57, 77–93.
- 724 Yamashita, R., Nishio, M., Do, R.K.G., Togashi, K., 2018. Convolutional
725 neural networks: an overview and application in radiology. *Insights
726 into imaging* 9, 611–629.

727 Yin, X., Liu, Q., Pan, Y., Huang, X., Wu, J., Wang, X., 2021. Strength of
728 stacking technique of ensemble learning in rockburst prediction with
729 imbalanced data: Comparison of eight single and ensemble models.
730 Natural Resources Research 30, 1795–1815.

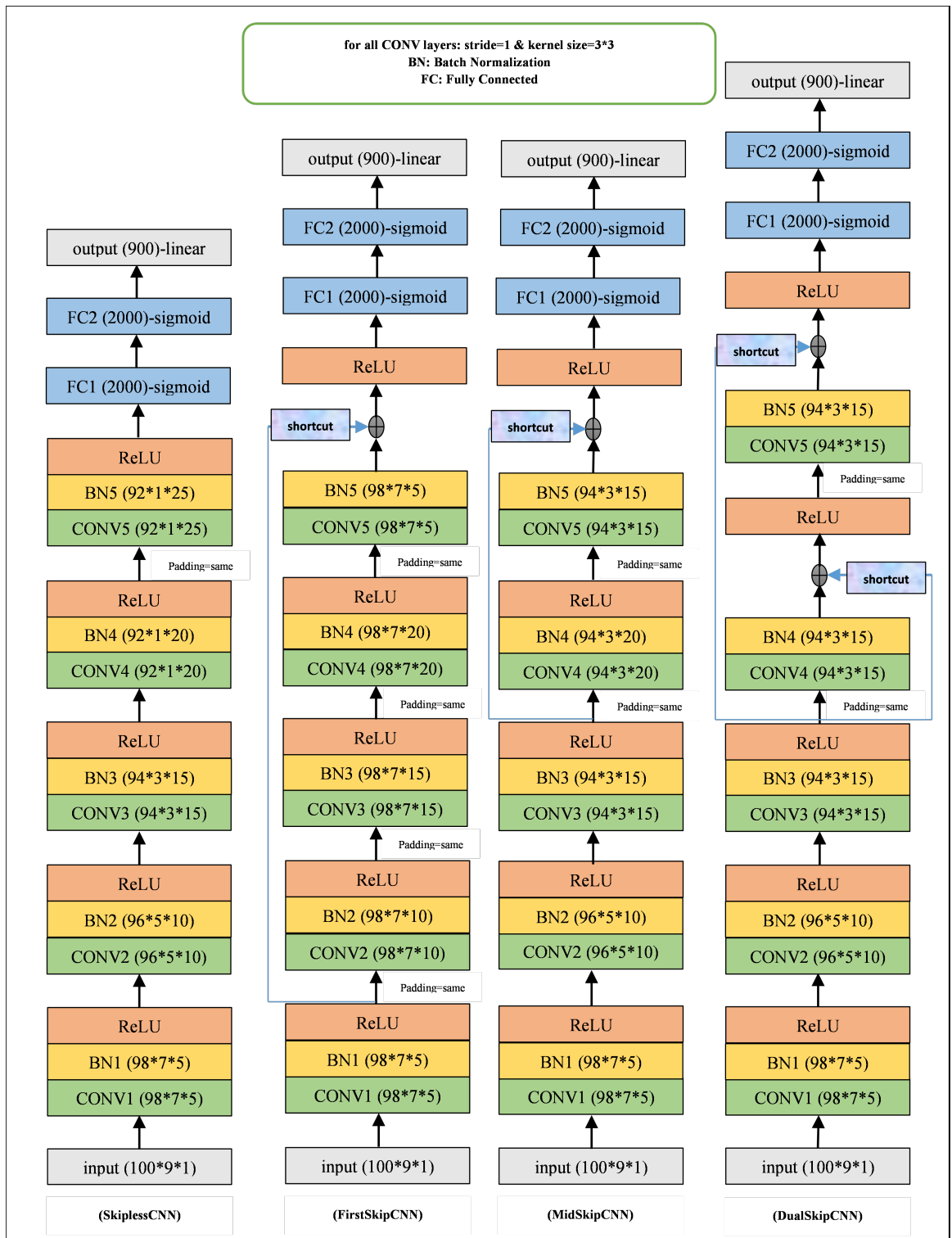


Figure 3: Structure of the CNN base learners configured in this study.

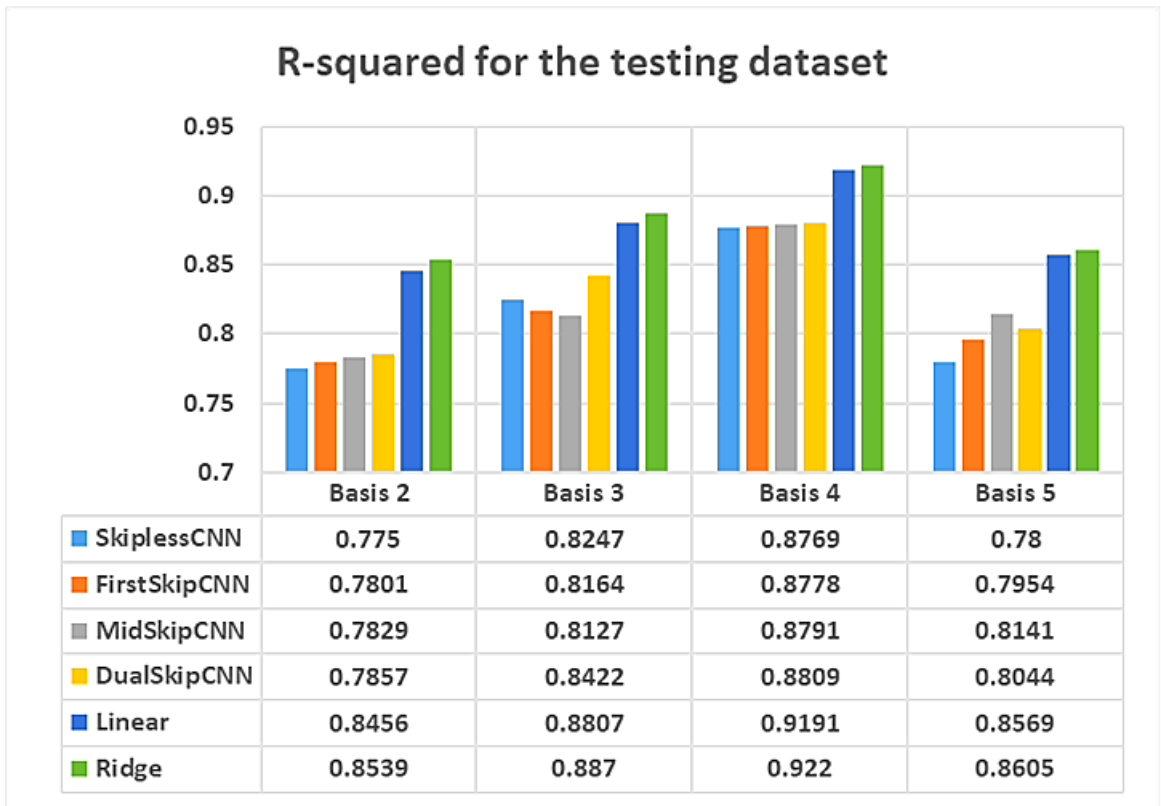


Figure 4: Prediction error analysis of the DEL-based and CNN models applied to the testing data subset, expressed in terms of R^2 .

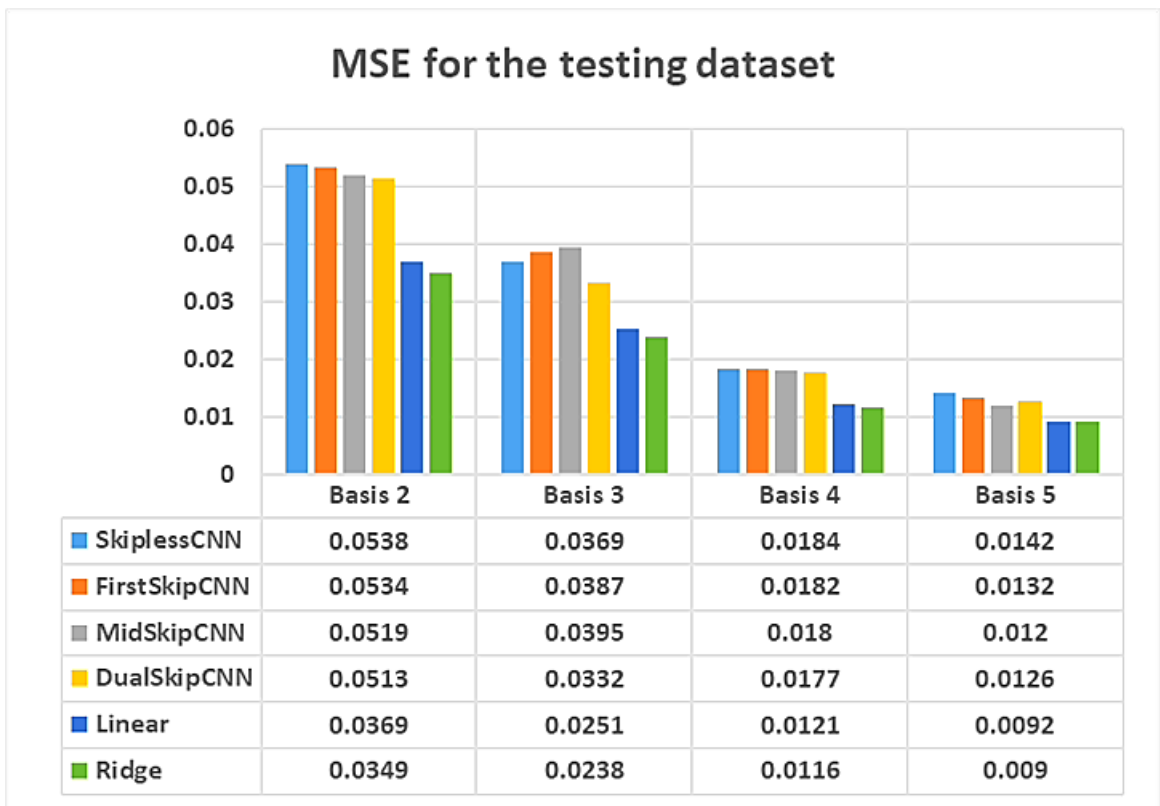


Figure 5: Prediction error analysis of the DEL-based and CNN models applied to the testing data subset, expressed in terms of *MSE*.

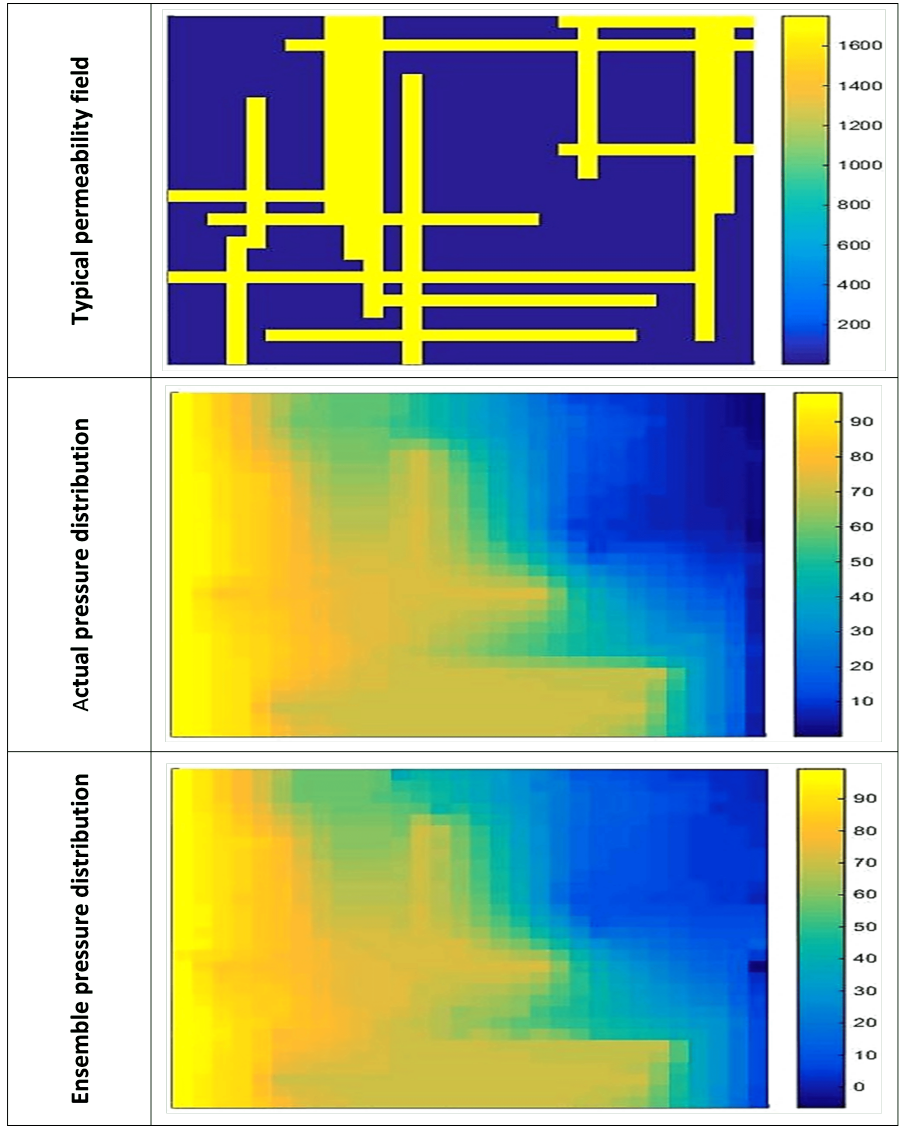


Figure 6: A comparison between the actual pressure distribution and the one obtained by ridge regression ensemble model for a representative permeability field.