

Machine Learning Based Localization for Iterative Ensemble Smoothers

Vinicius Silva, Gabriel Serrão, Alexandre Emerick

20th EnKF Workshop, Norway 2025



Localization

- Distance-based localization is highly effective. However, key reservoir parameters lack a direct spatial relationship with observations.
- ES-MDA with Schur-product localization of the Kalman gain:

$$\mathbf{m}_j^{\ell+1} = \mathbf{m}_j^\ell + \mathbf{R} \circ \underbrace{\left[\tilde{\mathbf{C}}_{\text{md}}^\ell \left(\tilde{\mathbf{C}}_{\mathbf{d}\mathbf{d}}^\ell + \alpha_\ell \mathbf{C}_{\mathbf{e}} \right)^{-1} \right]}_{\tilde{\mathbf{K}}^\ell} (\mathbf{d}_{\text{obs},j}^\ell - \mathbf{g}(\mathbf{m}_j^\ell))$$

- Several distance-free methods have been proposed for scalar parameters, although none of them is completely effective.

Pseudo-Optimal Localization*

- Minimization (term-by-term) of the Frobenius norm of the difference between the true covariance and its localized estimate

$$\mathbb{E} \left[\| \mathbf{C} - \mathbf{R} \circ \tilde{\mathbf{C}} \|^2_F \right] = \text{tr} (\mathbf{C}^2) - 2\mathbb{E} \left[\text{tr} \left(\mathbf{C} (\mathbf{R} \circ \tilde{\mathbf{C}}) \right) \right] + \mathbb{E} \left[(\mathbf{R} \circ \tilde{\mathbf{C}})^2 \right]$$

$$r_{ik} = \frac{c_{ik}^2}{c_{ik}^2 + \frac{c_{ik}^2 + c_{ii}c_{kk}}{N_e}} \quad \rightarrow \quad r_{ik} = \frac{\rho_{ik}^2}{\rho_{ik}^2 + \frac{1+\rho_{ik}^2}{N_e}}$$

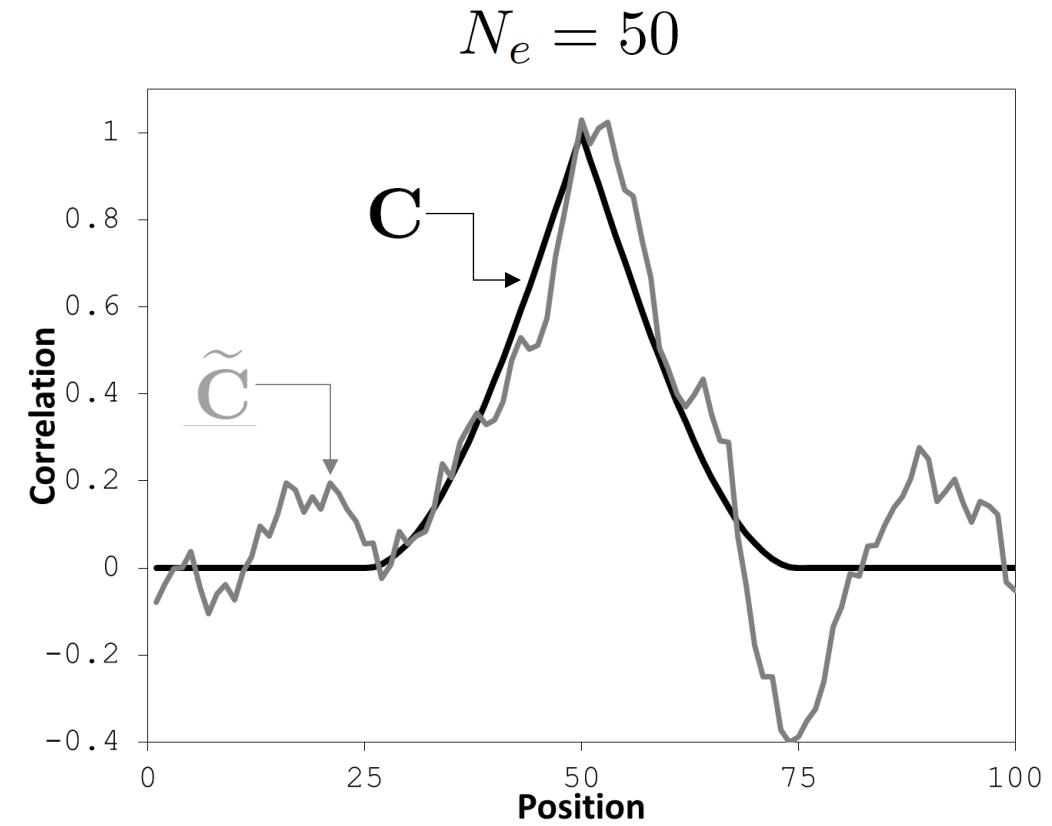
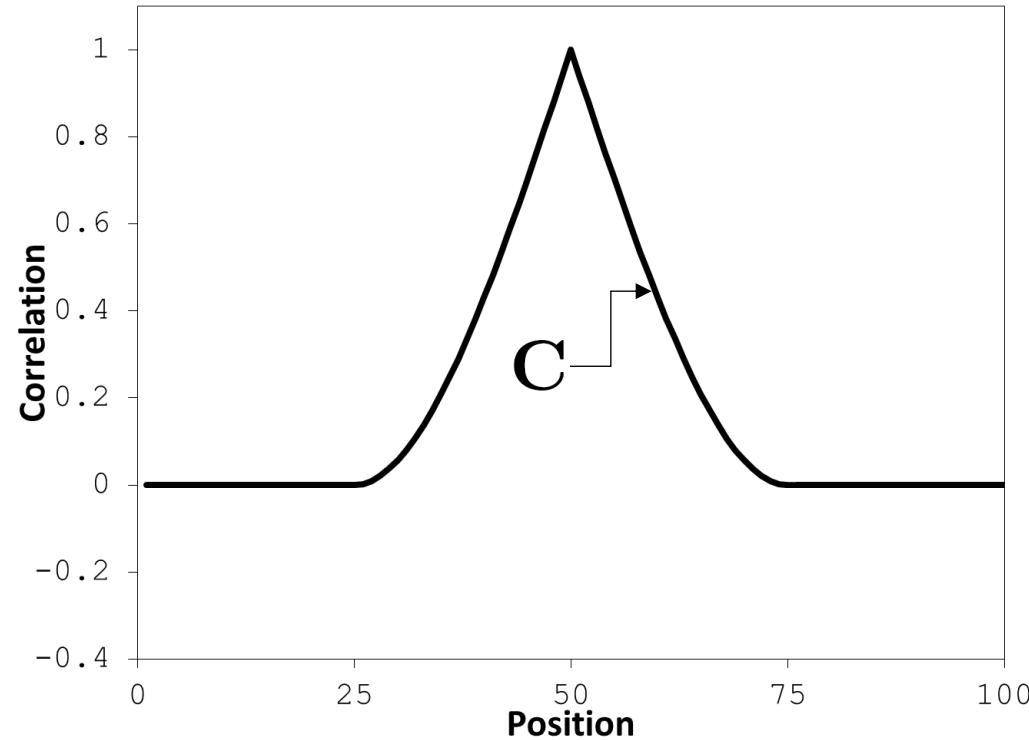
→

 $\begin{cases} \rho_{ik} & \text{correlation coefficient} \\ i \text{th} & \text{model parameter} \\ k \text{th} & \text{datum} \end{cases}$

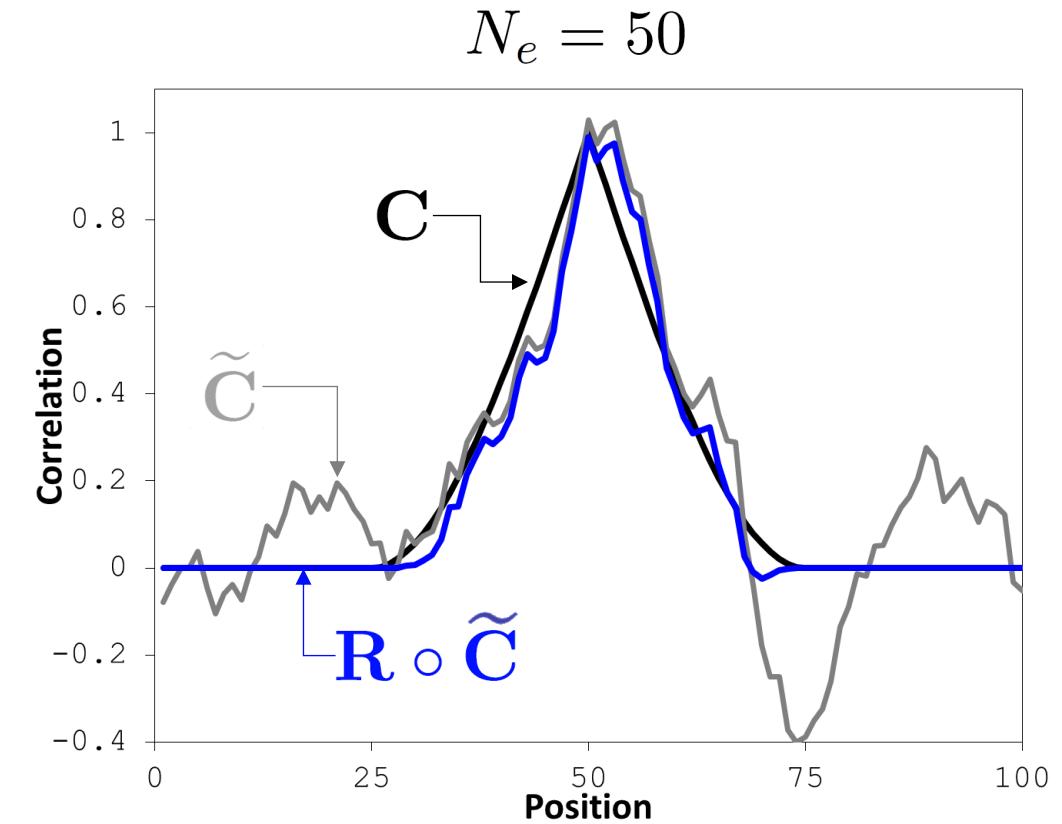
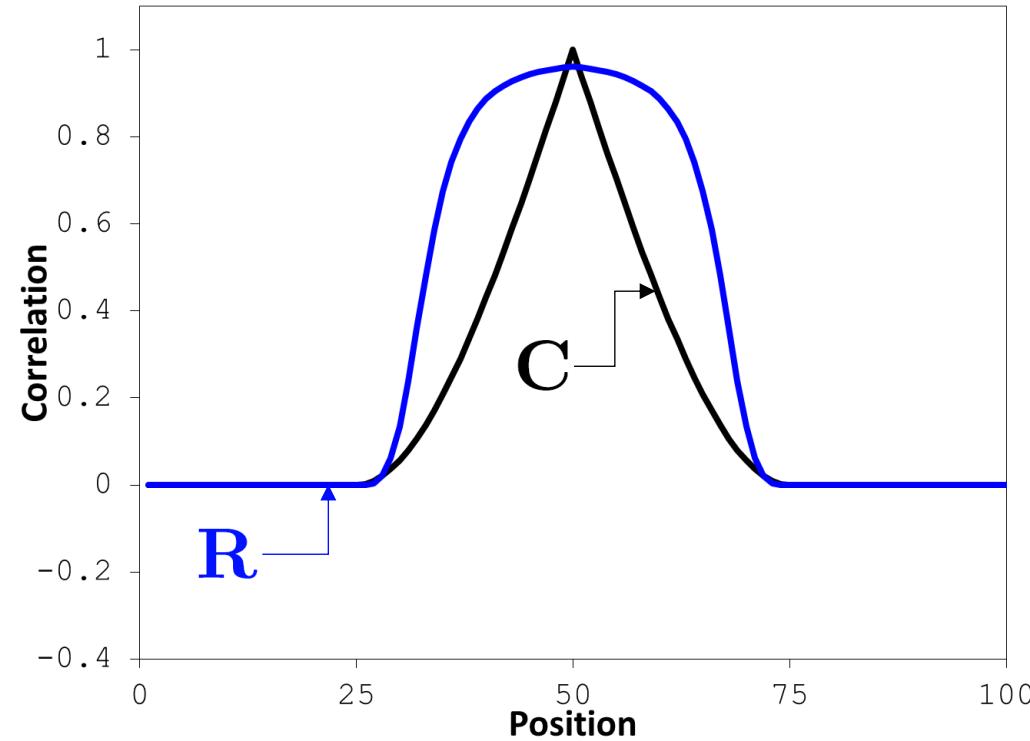
$$\tilde{r}_{ik} = \frac{\tilde{\rho}_{ik}^2}{\tilde{\rho}_{ik}^2 + \frac{1+\tilde{\rho}_{ik}^2}{N_e}} \quad \tilde{r}_{ik} = 0, \quad \text{if } |\tilde{\rho}_{ik}| < \tau$$

*Furrer, R., & Bengtsson, T.: Estimation of high-dimensional prior and posterior covariance matrices in Kalman filter variants, *Journal of Multivariate Analysis*, 98(2), 2007

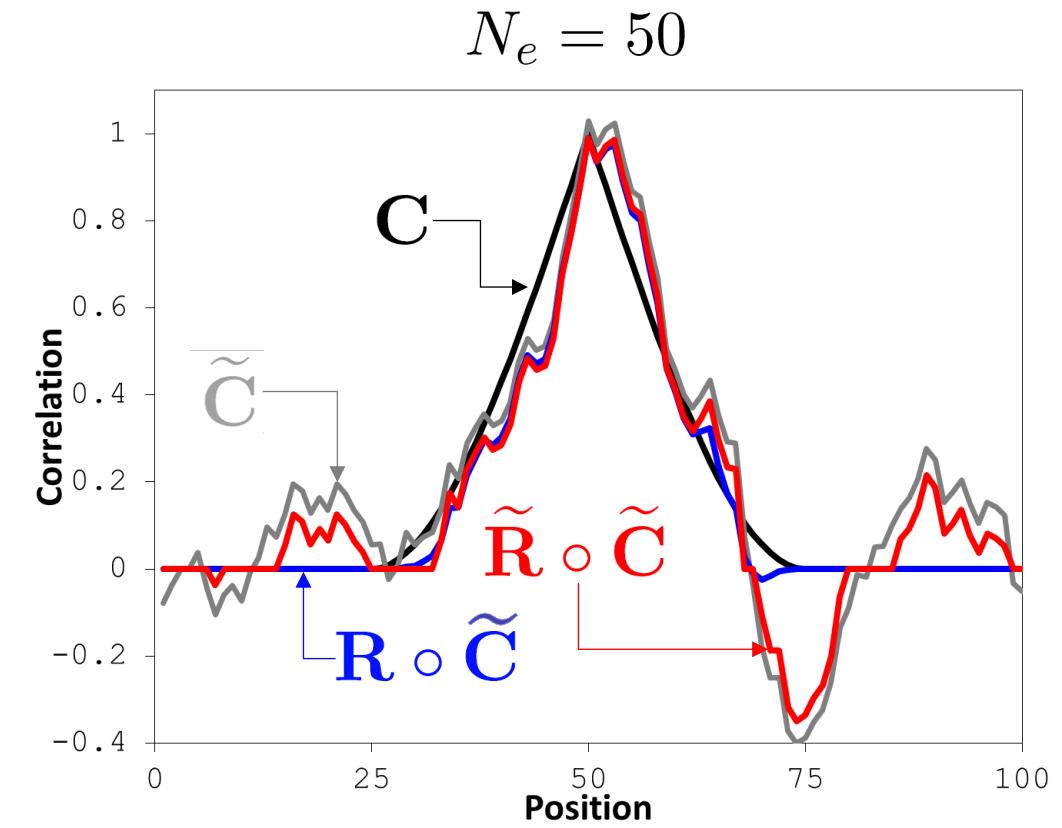
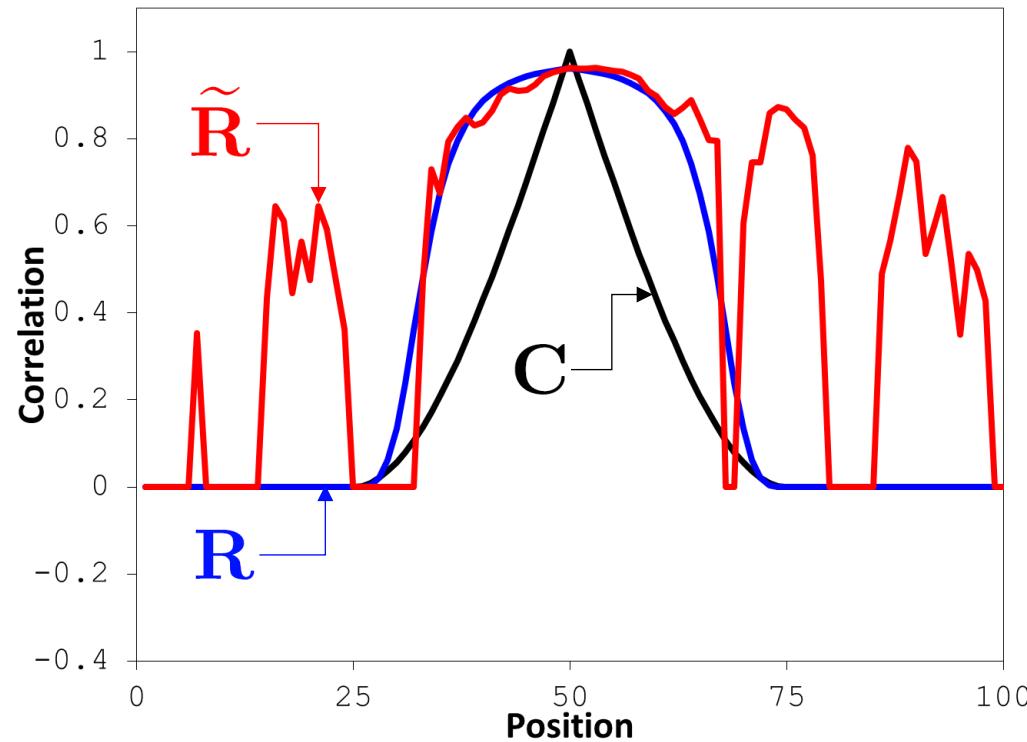
Pseudo-Optimal Localization



Pseudo-Optimal Localization



Pseudo-Optimal Localization

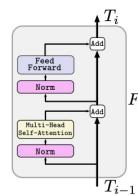
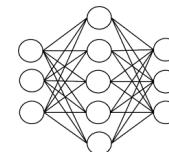
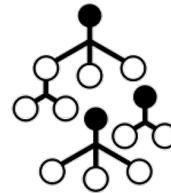


Machine Learning Localization

- 1 Simulate the prior ensemble: $\left\{ \mathbf{m}_j \right\}_{j=1}^{N_e} \xrightarrow{\text{Forward simulator}} \left\{ \mathbf{d}_j = \mathbf{g}(\mathbf{m}_j) \right\}_{j=1}^{N_e}$
- 2 Train the machine learning proxy: $\hat{\mathbf{d}} = \mathcal{G}(\mathbf{m})$
- 3 Use the ML proxy in a large ensemble: $\left\{ \mathbf{m}_j \right\}_{j=1}^{N_E} \xrightarrow{\text{ML proxy}} \left\{ \hat{\mathbf{d}}_j = \mathcal{G}(\mathbf{m}_j) \right\}_{j=1}^{N_E} \quad N_E \gg N_e$
- 4 Compute the cross-covariance: $\tilde{\mathbf{C}}_{\mathbf{m}\mathbf{d}} = \frac{1}{N_E - 1} \sum_{j=1}^{N_E} (\mathbf{m}_j - \bar{\mathbf{m}})(\hat{\mathbf{d}}_j - \bar{\hat{\mathbf{d}}})^\top$
- 5 Compute the localization coefficients: $\tilde{r}_{md} = \frac{\tilde{\rho}_{md}^2}{\tilde{\rho}_{md}^2 + \frac{1+\tilde{\rho}_{md}^2}{N_e}} \quad \tilde{\rho}_{md} = \frac{\tilde{c}_{md}}{\sqrt{\tilde{c}_{mm}\tilde{c}_{dd}}}$
 $\tilde{r}_{md} = 0, \quad \text{if } |\tilde{\rho}_{md}| < \tau = 10^{-3}$

Machine Learning Methods

- Support Vector Regression
- Extra Tree
- Decision Tree
- Random Forest
- XGBoost
- LightGBM
- Neural Network (Multilayer Perceptron)
- Tabnet
- FT-Transformer



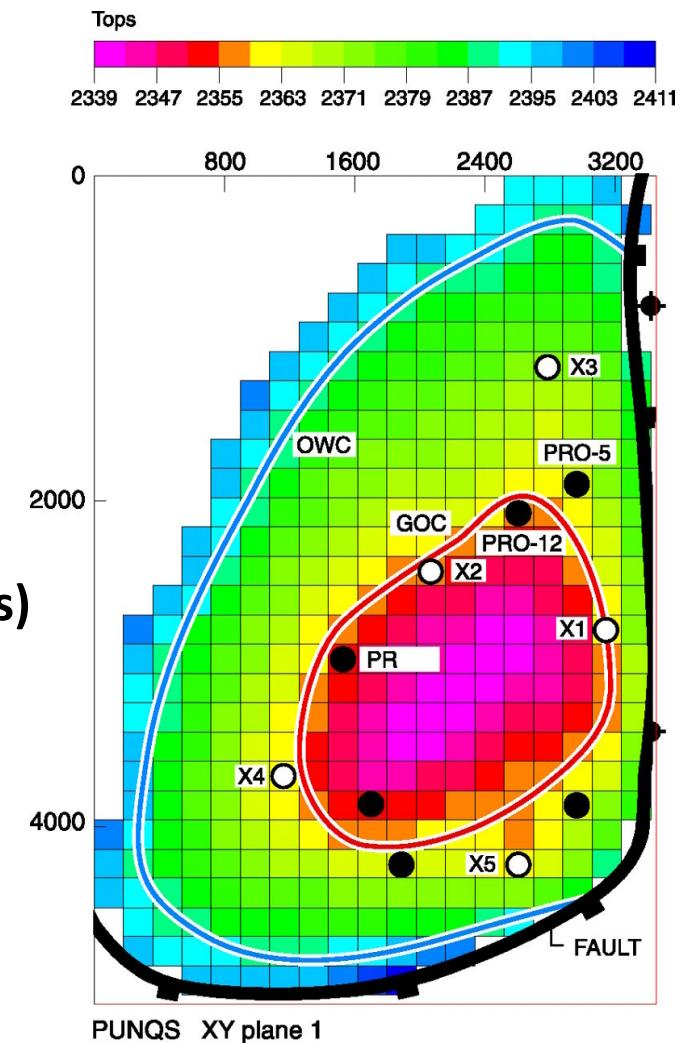
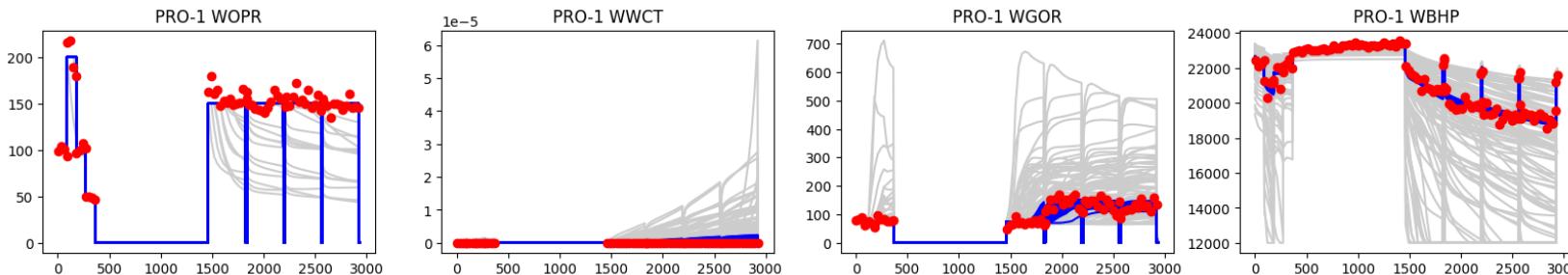
Goals:

- No additional reservoir simulations (beyond those already used in the data assimilation)
- No additional tuning parameter (no hyperparameter tuning)

Part 1 – Localization for Scalar Parameters

Test Case (PUNQ-S3* modified model)

- Model parameters (20 scalar variables)
 - 5 Porosity multipliers
 - 5 Permeability multipliers
 - Compressibility
 - Water-oil contact depth
 - Gas-oil contact depth
 - 2 Aquifers radius
 - **5 Dummy variables**
- Observed data from 6 producers -- 8 years history (1530 observed points)



*Floris, F.J., Bush, M.D., Cuypers, M., Roggero, F., & Syversveen, A.R.: Methods for Quantifying the Uncertainty of Production Forecasts: a Comparative Study, *Petroleum Geoscience*, 7(S), 2001

Results

- Ensemble size (training set): $N_e = 100$
- Reference: $N_E = 5000$

$$\text{RMSE} = \sqrt{\frac{1}{N_m N_d} \sum_{i=1}^{N_m} \sum_{k=1}^{N_d} (\tilde{\rho}_{ik}^{5000} - \tilde{r}_{ik} \cdot \tilde{\rho}_{ik}^{100})^2}$$

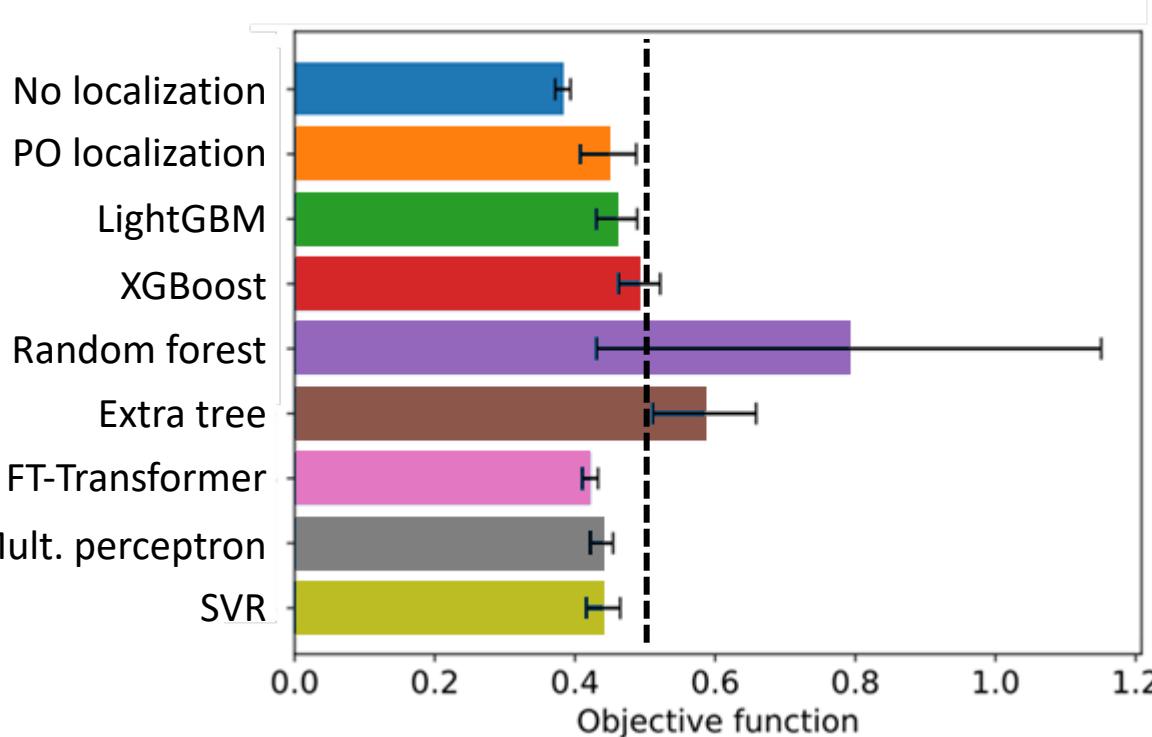
	Training time (s)	RMSE (Correlation)
No localization	-	0.053
Pseudo-optimal localization	-	0.047
Decision Tree	0.2	0.052
Support Vector Regression	2.7	0.041
Random Forest	10.3	0.046
Extra Tree	6.3	0.044
XGBoost	58.1	0.043
LightGBM	31.7	0.040
Multilayer Perceptron	111	0.039
FT-Transform	269.8	0.039
Tabnet	50.9	0.051

Data Assimilation Results

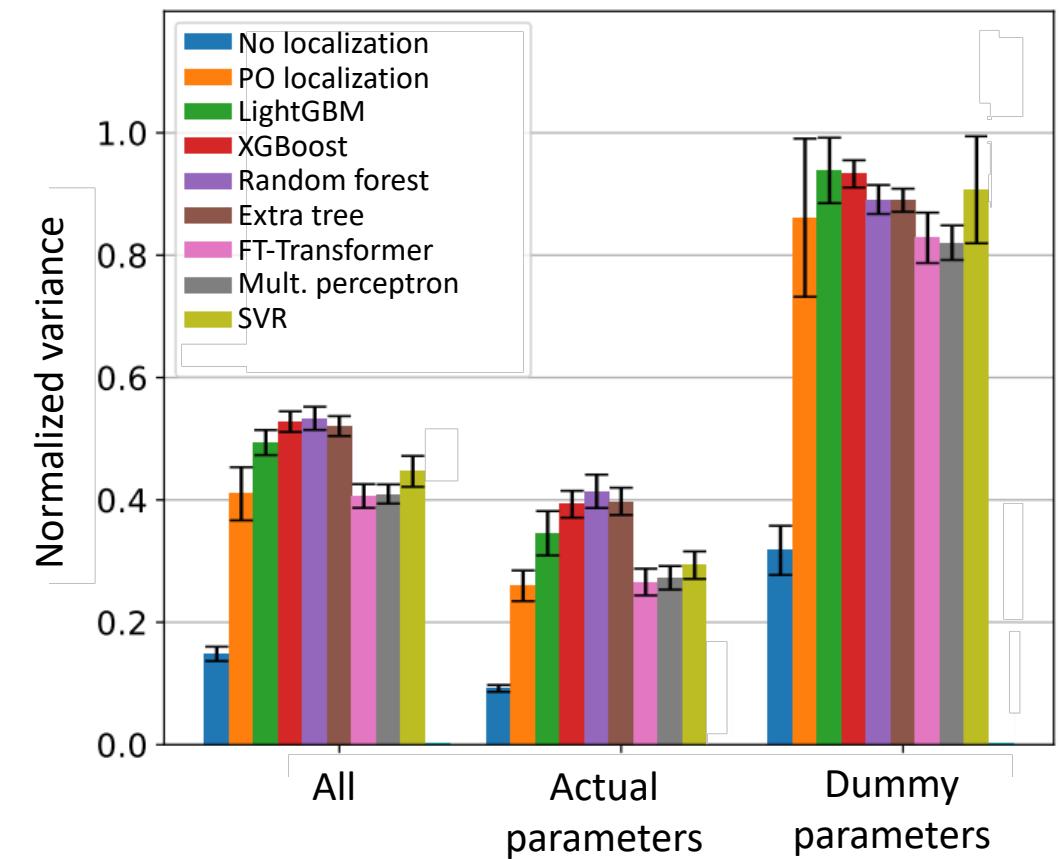
- 10 data assimilation experiments ($N_e = 100$)

$$NV = \frac{\text{var} [\mathbf{m}_{\text{post}}]}{\text{var} [\mathbf{m}_{\text{prior}}]}$$

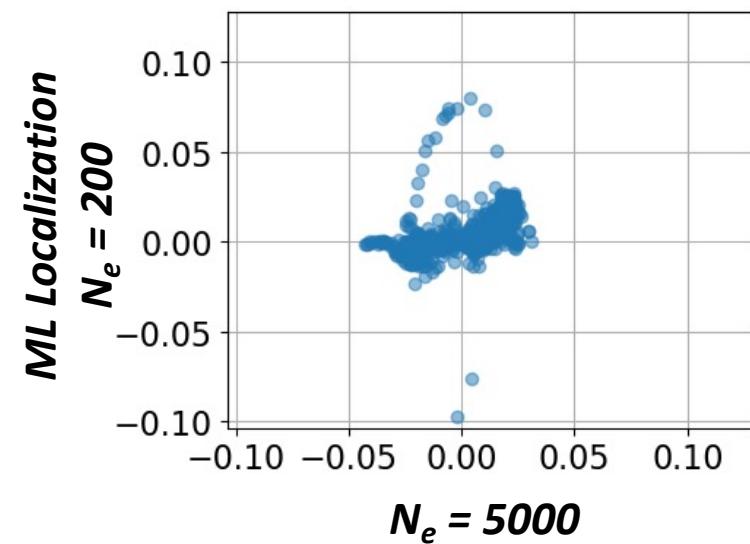
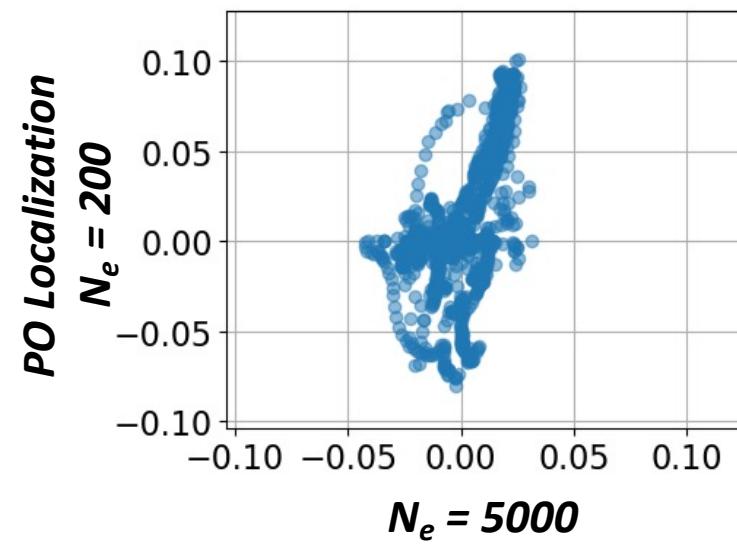
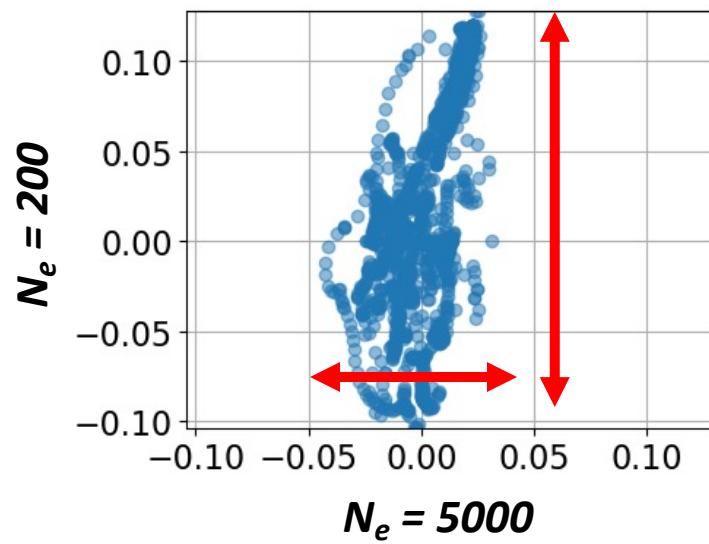
Data Mismatch



Normalized Variance



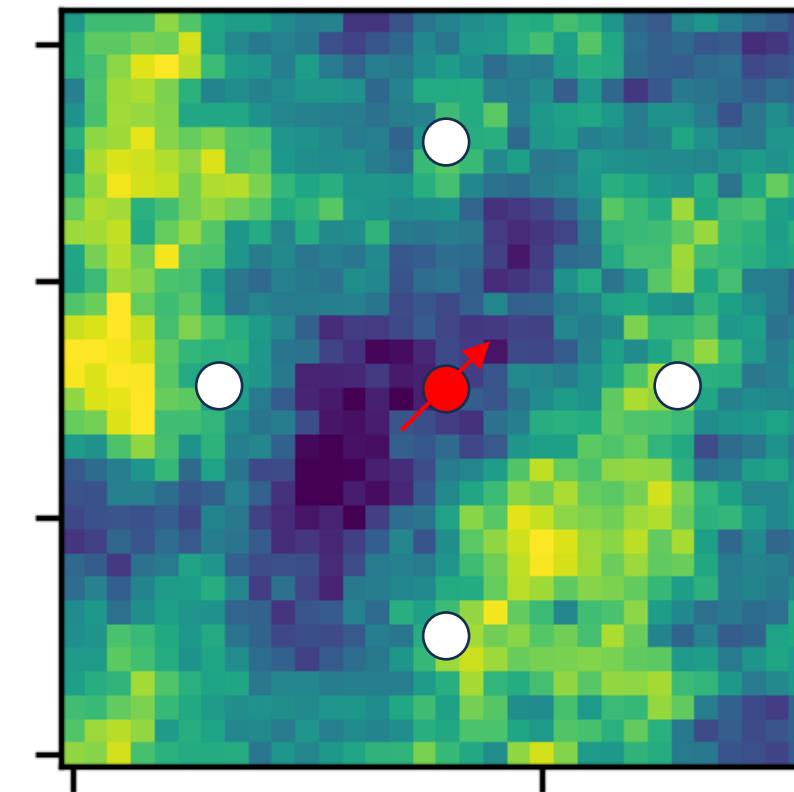
Correlations (Dummy vs Predicted Data)



Part 2 – Localization for Grid Parameters

Test Case 1

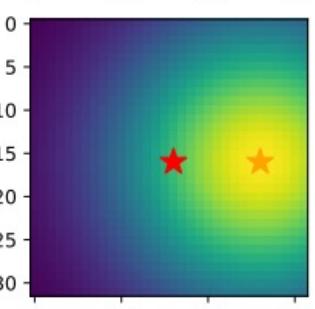
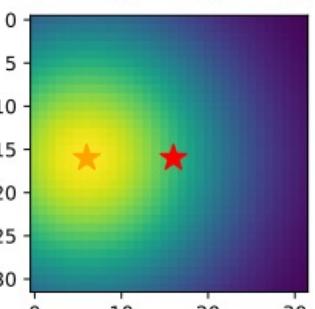
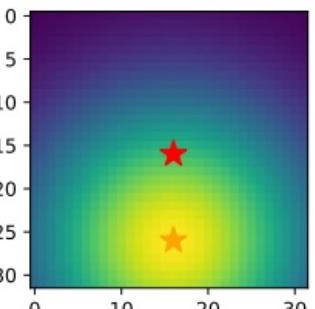
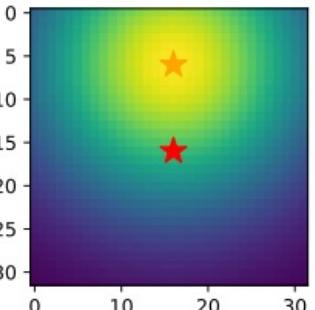
- CO₂ injection
- Pressure measurements at 4 monitor wells
- Estimate log-permeability distribution
- Ensemble with $N_e = 100$
- ML method: XGBoost



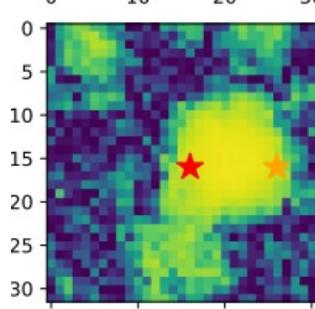
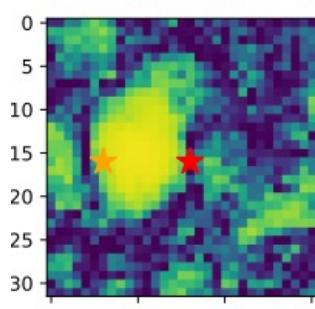
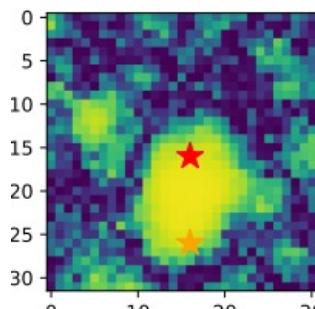
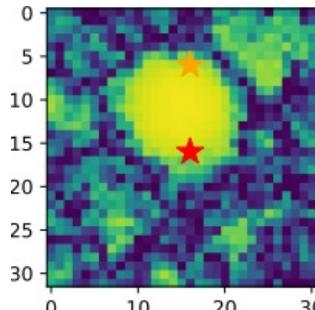
Localization Results



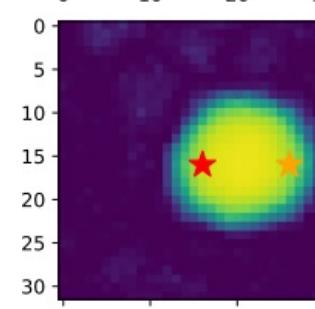
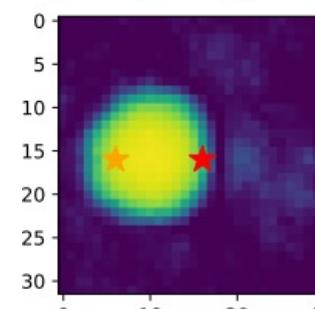
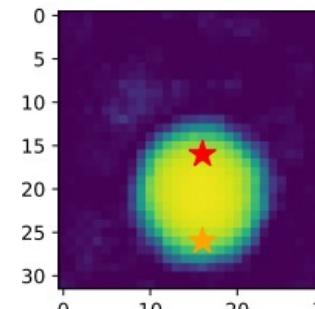
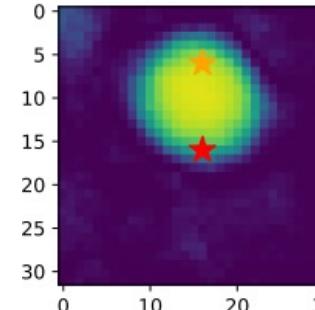
Distance-based localization



Pseudo-optimal localization



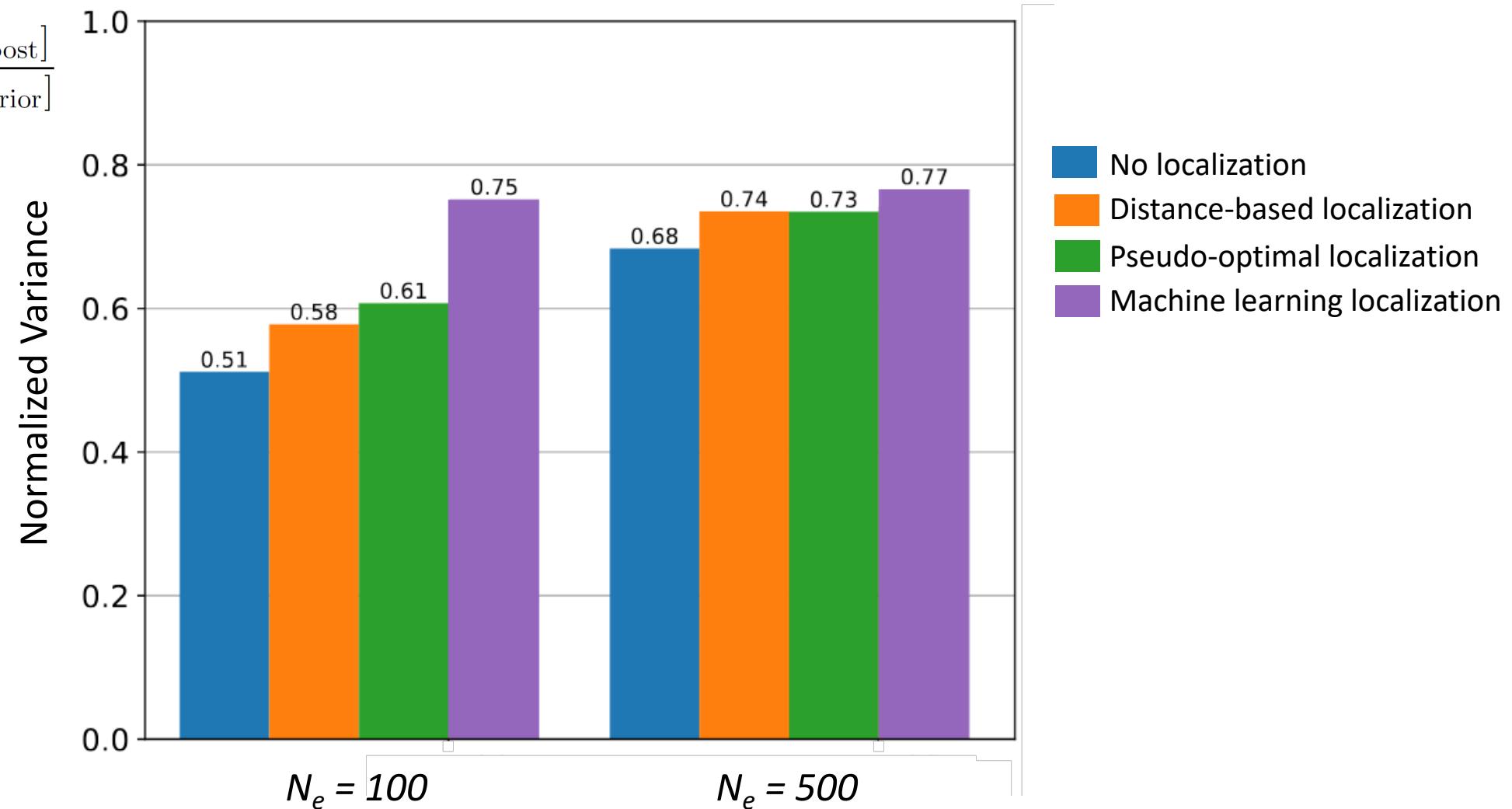
Machine learning localization



Data Assimilation Results

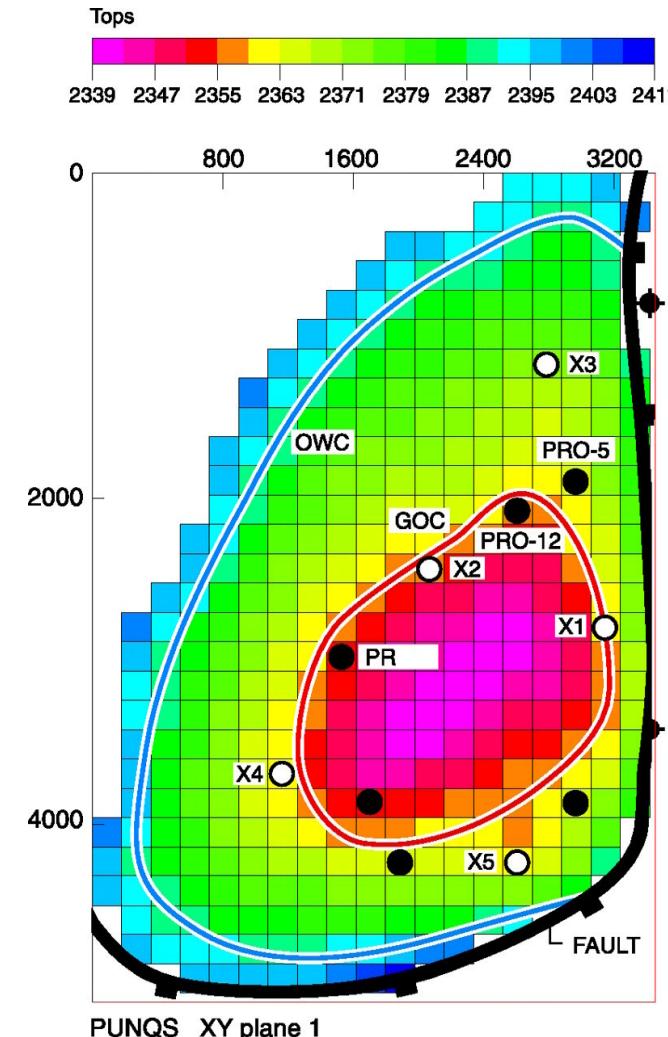


$$NV = \frac{\text{var} [\mathbf{m}_{\text{post}}]}{\text{var} [\mathbf{m}_{\text{prior}}]}$$

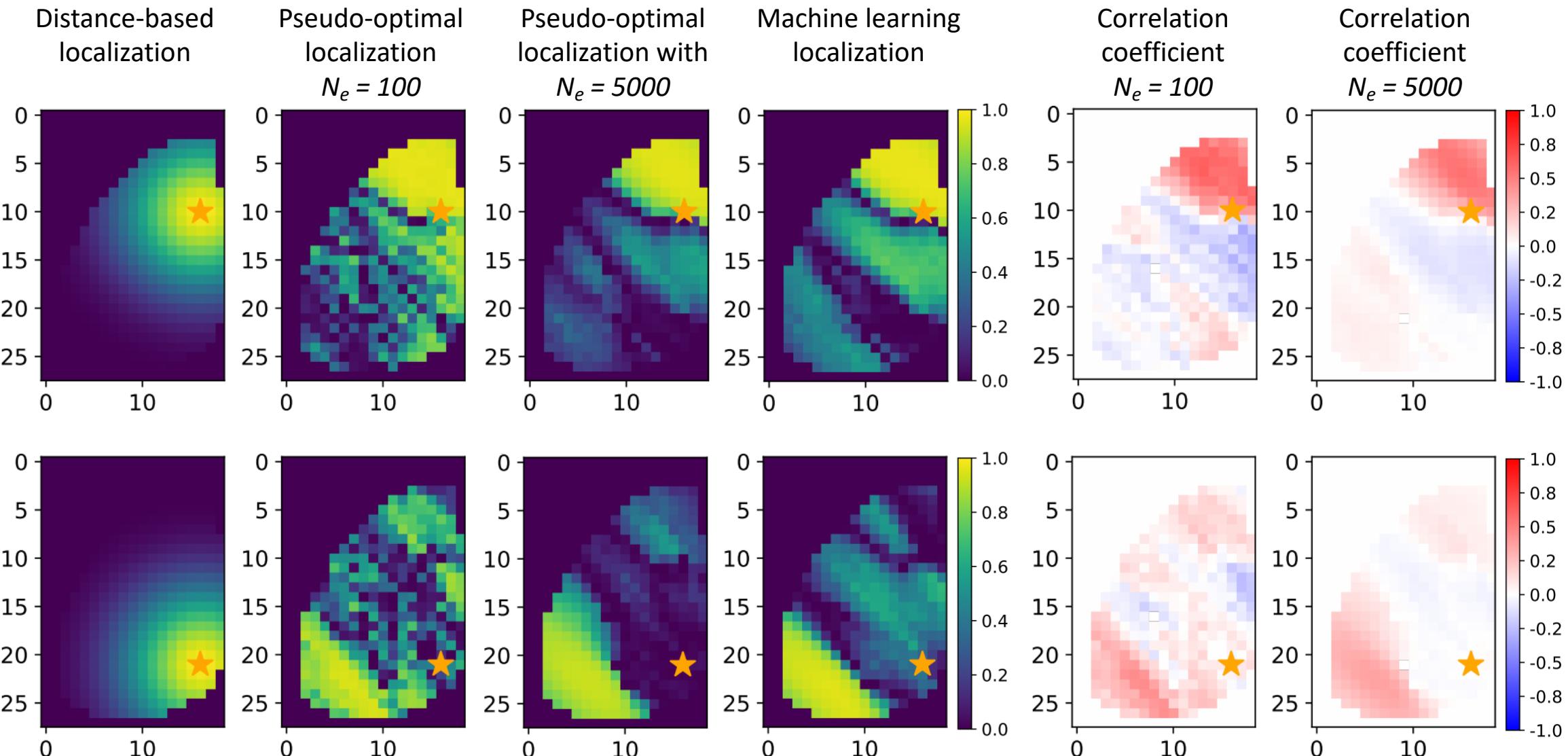


Test Case 2 (PUNQ-S3 modified model)

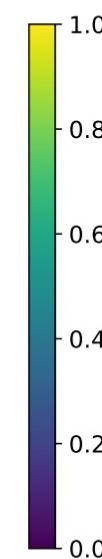
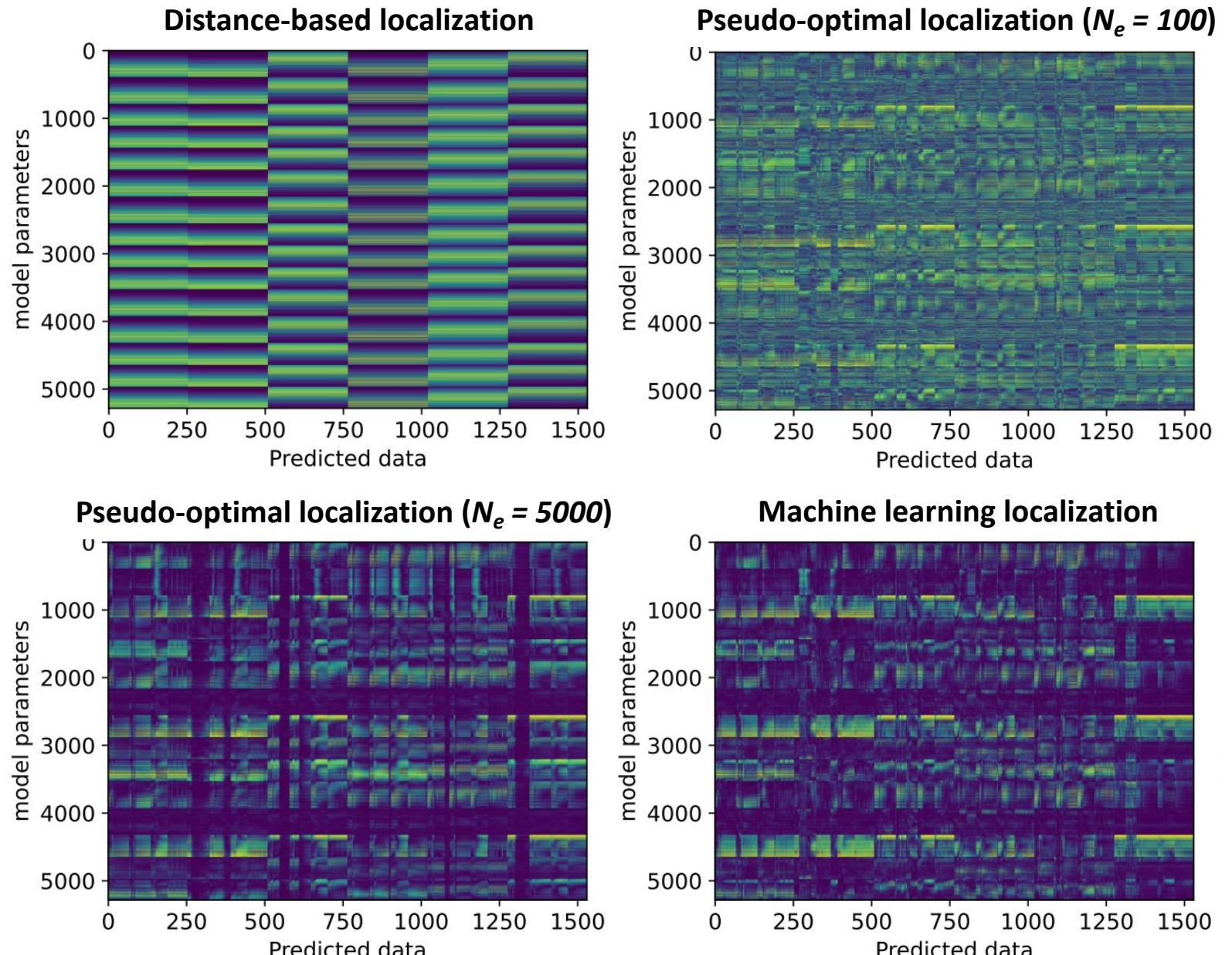
- Estimate porosity, horizontal and vertical permeability distribution
- Measurements of water cut, gas-oil ratio, and pressure
- Ensemble with $N_e = 100$
- ML method: LightGBM



Localization Results – Pressure

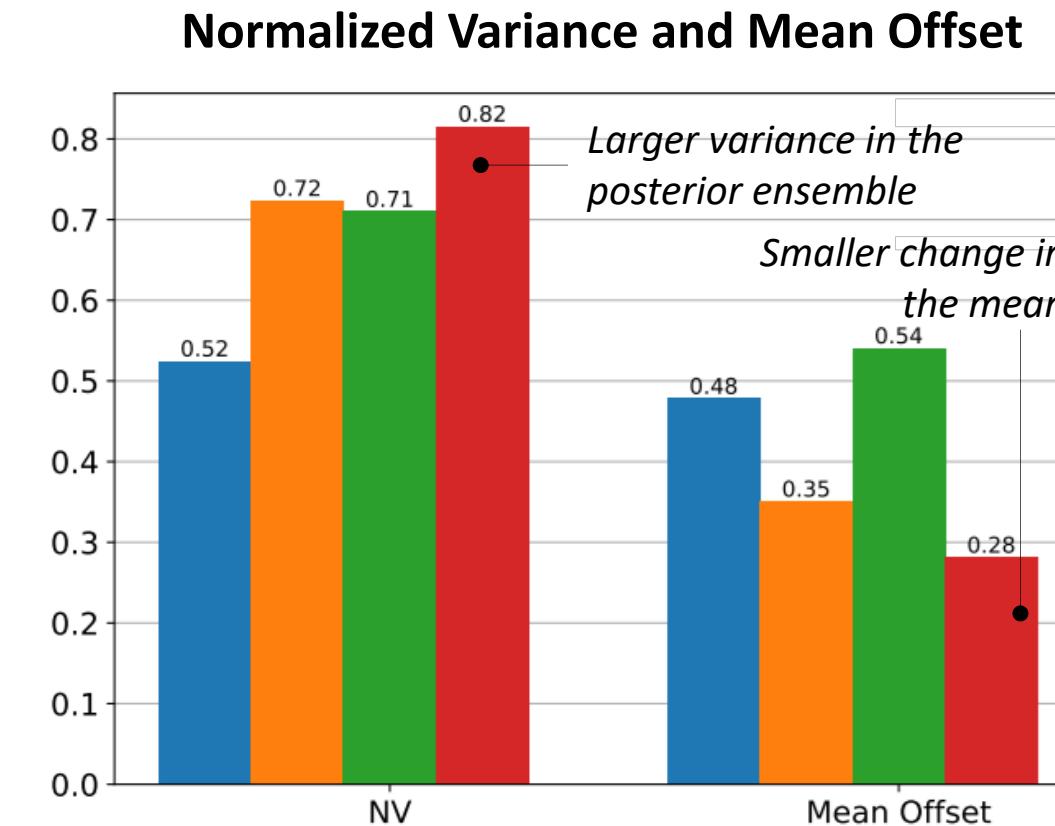
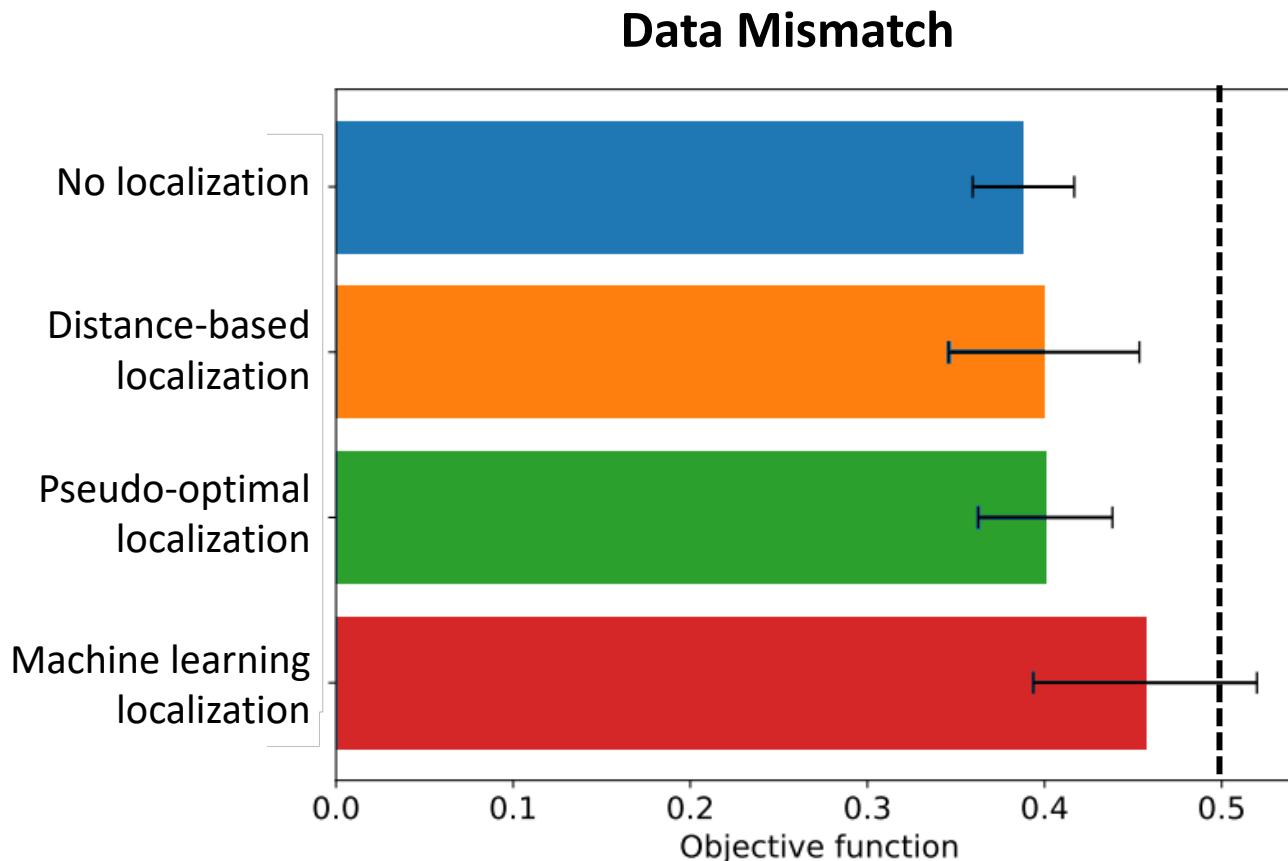


Localization Matrix



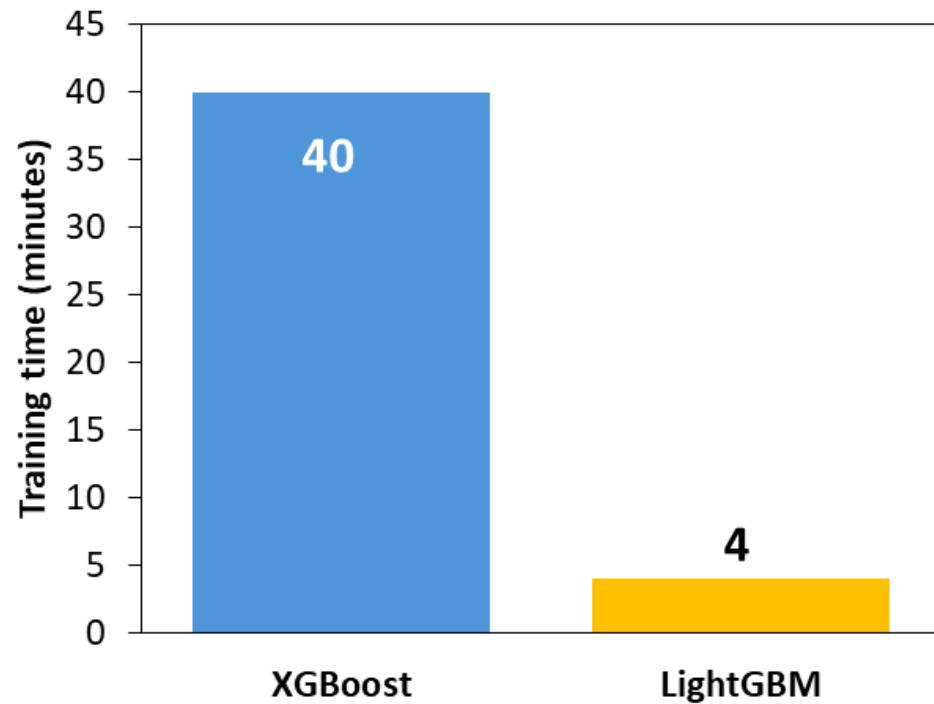
Data Assimilation Results

- 10 data assimilation experiments ($N_e = 100$)



Training Time*

- Number of model parameters: 5238
- Number of data points: 1530
- CPU Intel(R) Xeon(R) with 16 cores



*Preliminary result

Final Comments

- Distance-free localization scheme
- No additional forward simulations
- Out-of-the-box machine learning implementations (XGBoost or LightGBM)
- No additional tuning parameters
- Scalability for large scale problems needs further investigation

References

- Furrer, R., & Bengtsson, T.: **Estimation of high-dimensional prior and posterior covariance matrices in Kalman filter variants**, *Journal of Multivariate Analysis*, 98(2), 2007
- Floris, F.J., Bush, M.D., Cuypers, M., Roggero, F., & Syversveen, A.R.: **Methods for Quantifying the Uncertainty of Production Forecasts: a Comparative Study**, *Petroleum Geoscience*, 7(S), 2001
- Silva, V.L.S., Seabra, G.S. & Emerick, A.A.: **Machine Learning to Enhance the Covariance Estimations of Non-Local Model Parameters in Ensemble Based-Data Assimilation**, *Proceedings of the SPE Reservoir Simulation Conference*, SPE-223908-MS, 2025
- Silva, V.L.S., Seabra, G.S. & Emerick, A.A.: **Mitigating Loss of Variance in Ensemble Data Assimilation: Machine Learning-Based and Distance-Free Localizations for Better Covariance Estimation**, *In preparation*, 2020