# Improving Machine Learning Explainability with a Graphical User Interface

Cody McBroom Walker

Changing the World's Energy Future

**Idaho National Laboratory**

# Improving Machine Learning Explainability with a Graphical User Interface

**Cody McBroom Walker**

**March 2024**

**Idaho National Laboratory**
**Idaho Falls, Idaho 83415**

**http://www.inl.gov**

# Explainability comes in different forms depending on which model you use.

- How do we explain these models to the user?

- How much would you have to explain to go from an input to an output?



2 [m, c]



2 [2 splits]

| Data | Hyperparameters & Optimization | Model |
|------|-------------------------------|-------|

| Post hoc |
|----------|

Effects model performance and explainability

Explainability



Support vectors, kernel trick and hyperplane.



Number of weights, biases, & connections.

# LIME is a post-hoc method for black-box models.

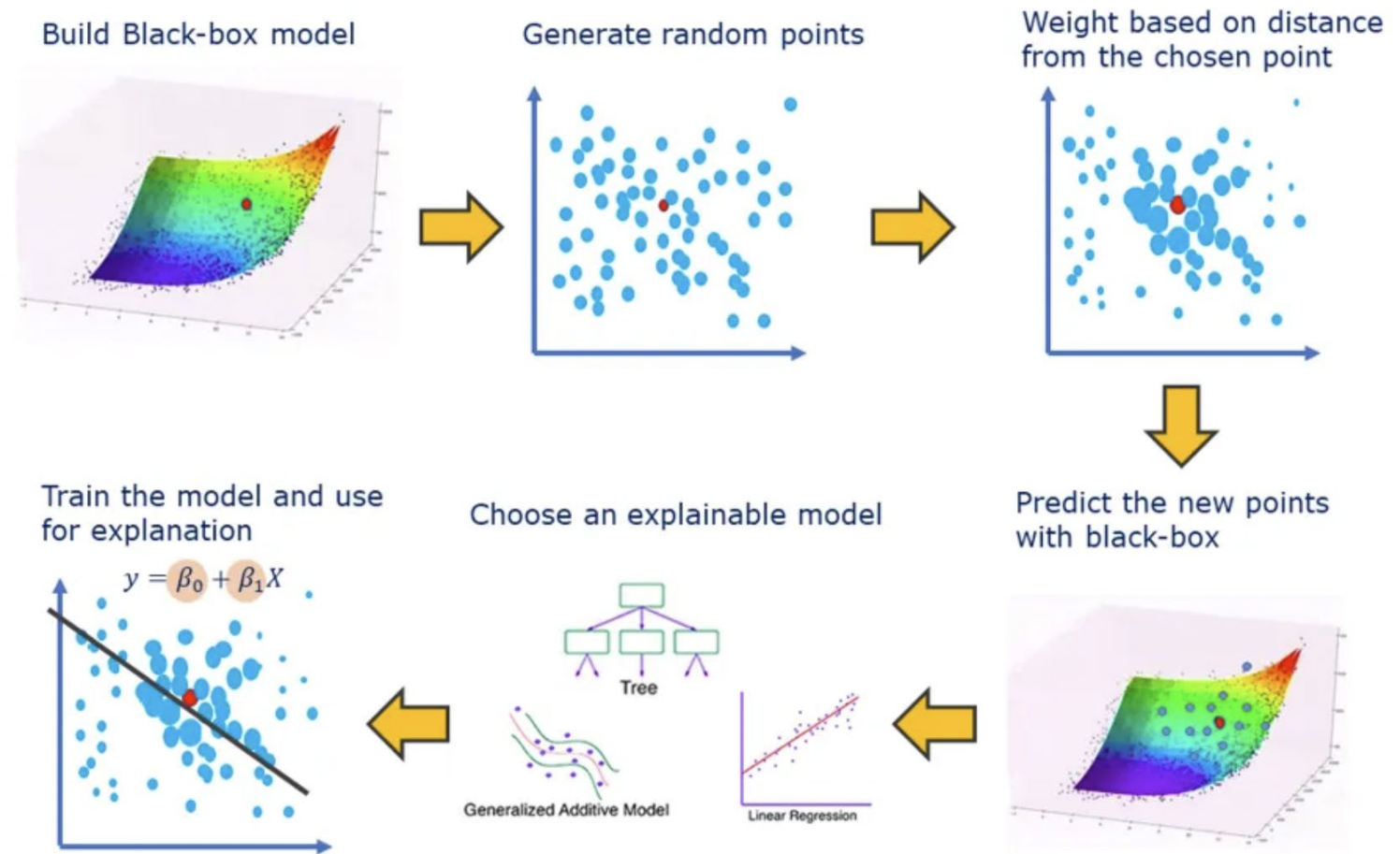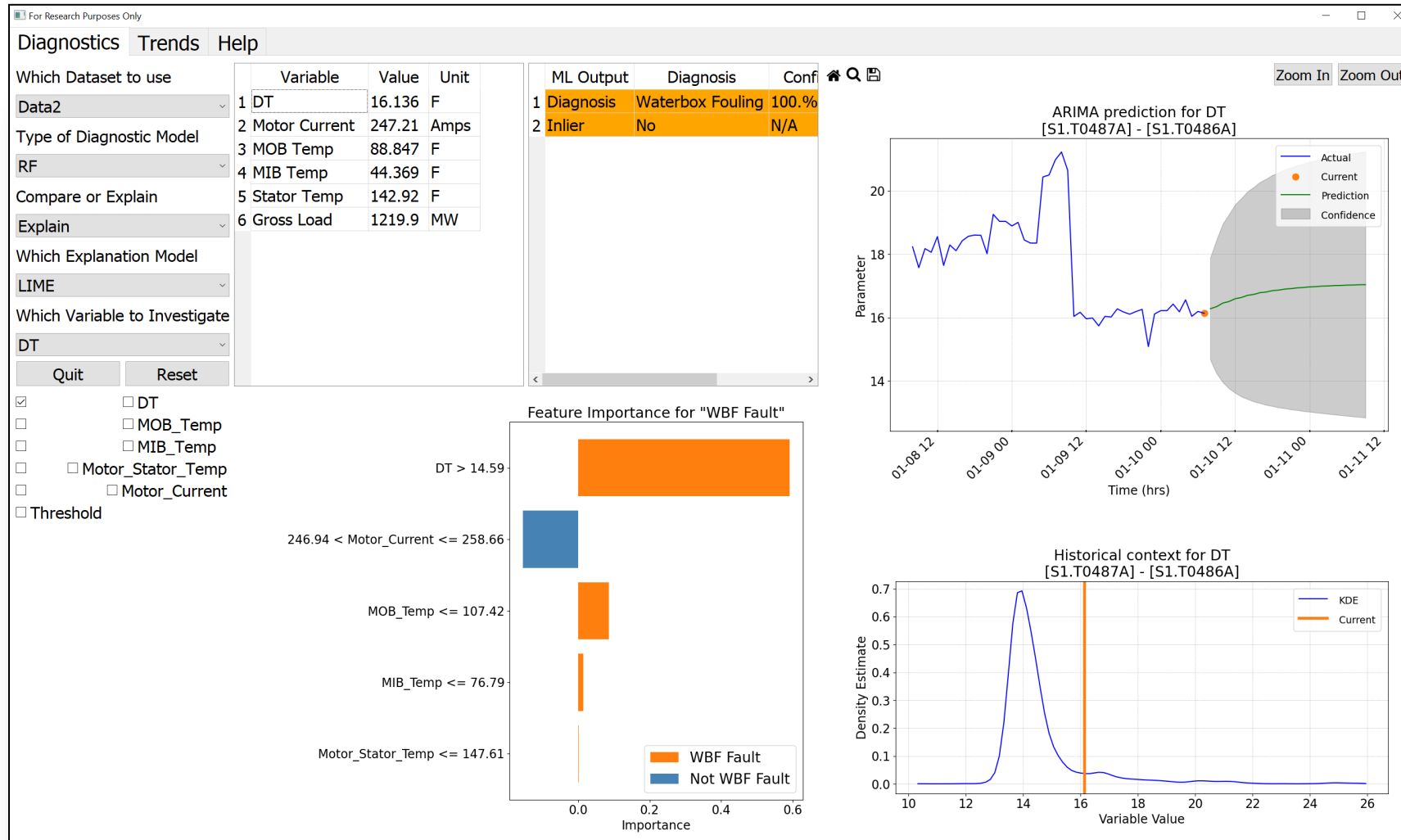- Local interpretable model-agnostic explanations (LIME) can be used for any model.

- LIME is only valid **locally**.

- SHAP (Shapley Additive Explanations) are another common post-hoc method used to increase explainability.



Build Black-box model

Generate random points

Weight based on distance from the chosen point

Predict the new points with black-box

Choose an explainable model

Tree

Generalized Additive Model

Linear Regression

Train the model and use for explanation

$$y = \beta_0 + \beta_1 X$$

Giorgio Visani, 2020 "LIME: explain Machine Learning predictions." Accessed 2024.
https://towardsdatascience.com/lime-explain-machine-learning-predictions-af8f18189bfe

# Model confidence, prognostics, explainability, and historical context all provide evidence for the conclusion.

# Adding context to the data can further improve understanding.