

What is explainable AI?

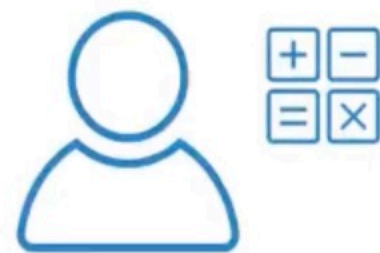


# **Interpretable Machine Learning**

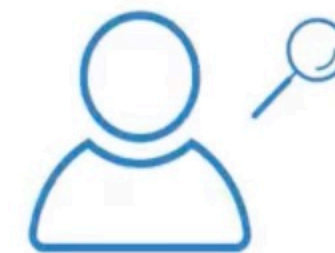
# What is explainable AI?



Source: Google Trends for "Explainable AI"

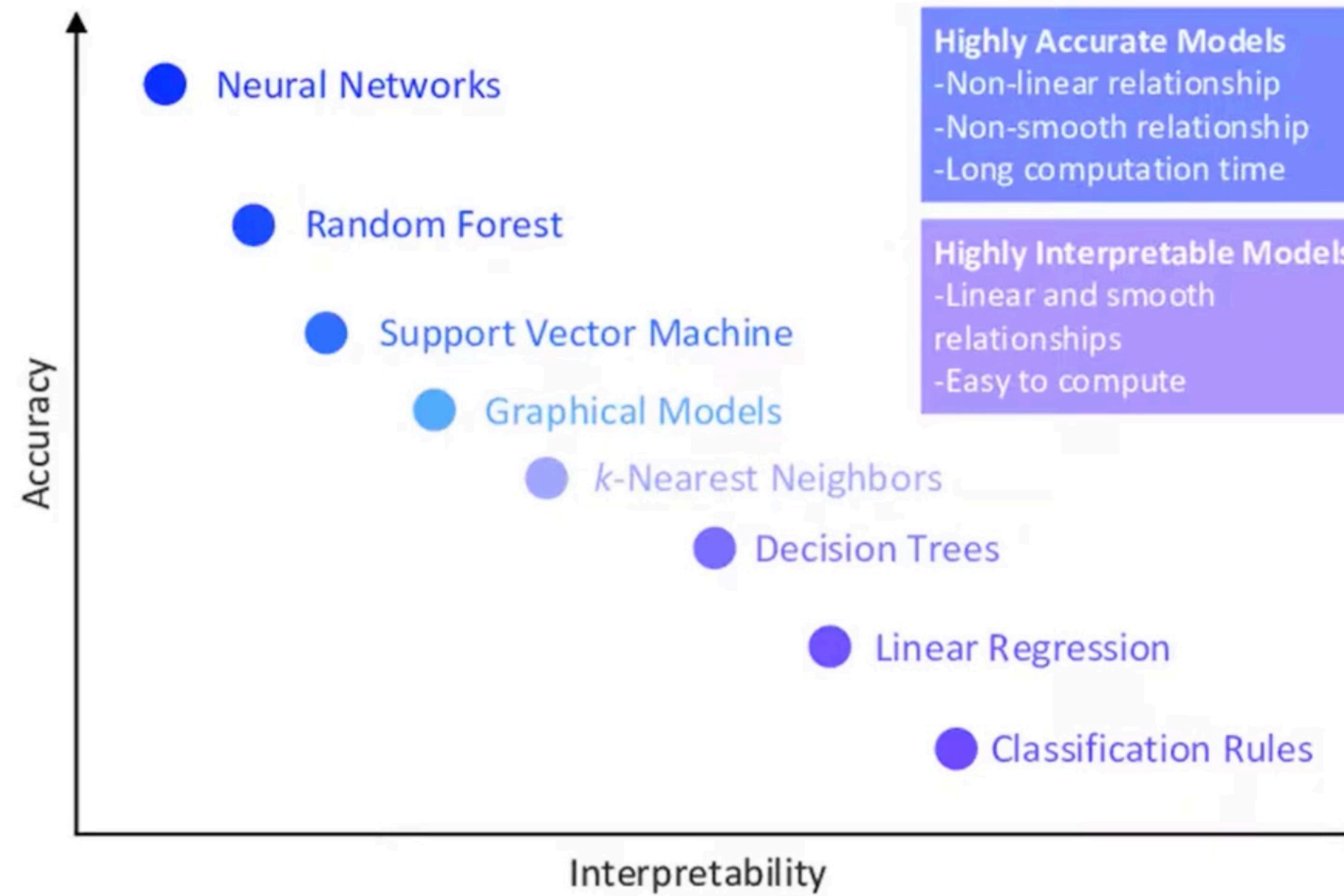


Data Scientists &  
ML Engineers



**Users** of the  
algorithms

# What is explainable AI?



Source: Machine Learning for 5G/B5G Mobile and Wireless Communications: Potential, Limitations, and Future Directions

# Understanding Machine Learning models



„Model based“

Build **interpretable** ML  
models



„Post-hoc“

Derive explanations for  
**complex** ML models

Black-box  
approach



White-box  
approach

# Categorization in XAI



**Model-agnostic**

Applicable to all model types

**Model-specific**

Only applicable to a specific model type

Agnosticity

**Global**  
explanation

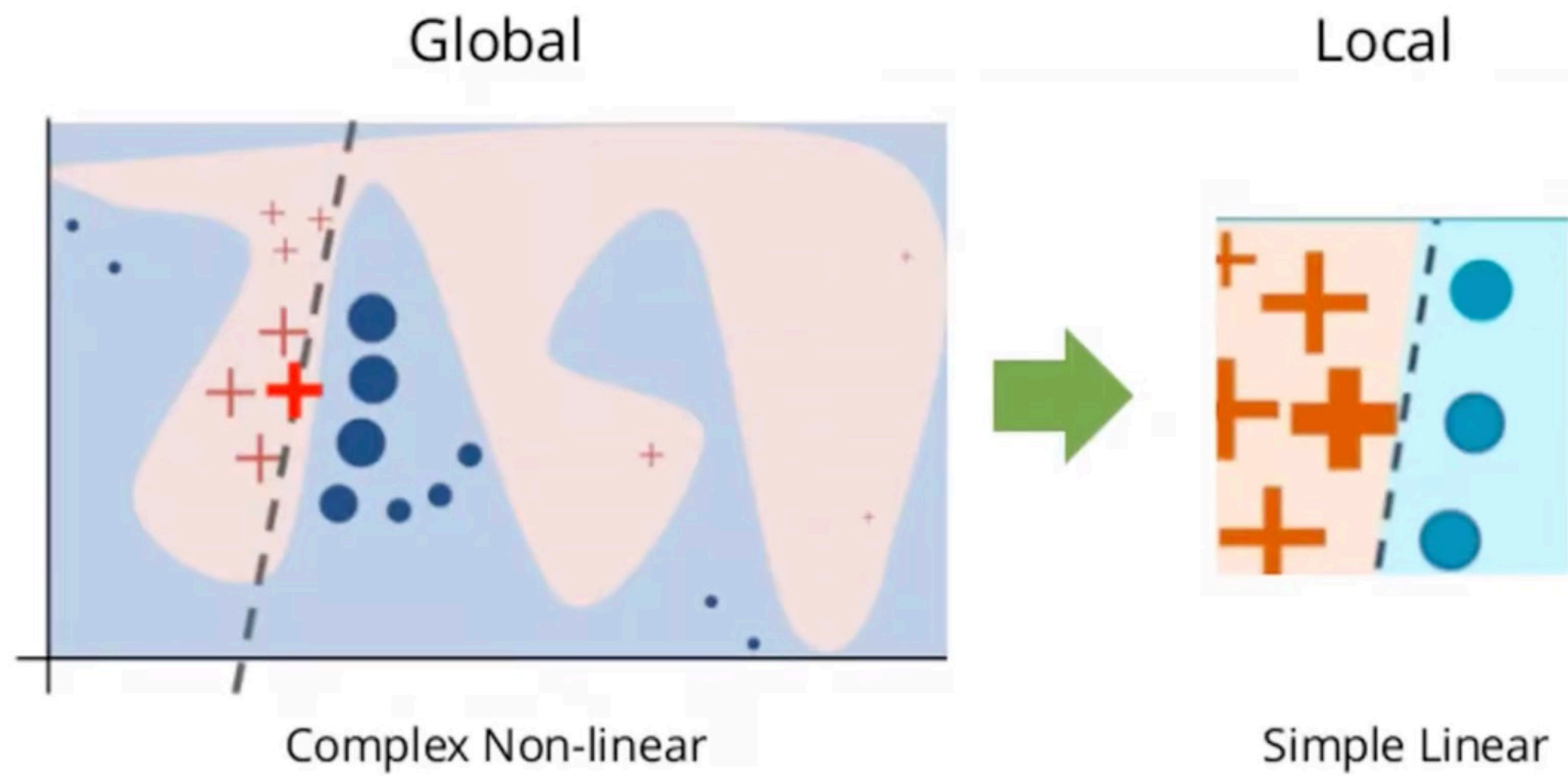
Explaining the whole model

**Local**  
explanation

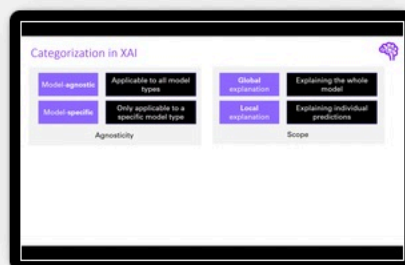
Explaining individual predictions

Scope





**Source:** "Why Should I Trust You?": Explaining the Predictions of Any Classifier, Ribeiro et al.



# Categorization in XAI



**Model-agnostic**

Applicable to all model types

**Model-specific**

Only applicable to a specific model type

Agnosticity

Graph

Image

Text / Speech

Tabular

Data type

**Global**  
explanation

Explaining the whole model

**Local**  
explanation

Explaining individual predictions

Scope



Visual

Feature  
importance

Data points

Surrogate  
models

Explanation type



# Explainable AI

Cheat sheet [ex.pegg.io](https://ex.pegg.io) v.0.2

