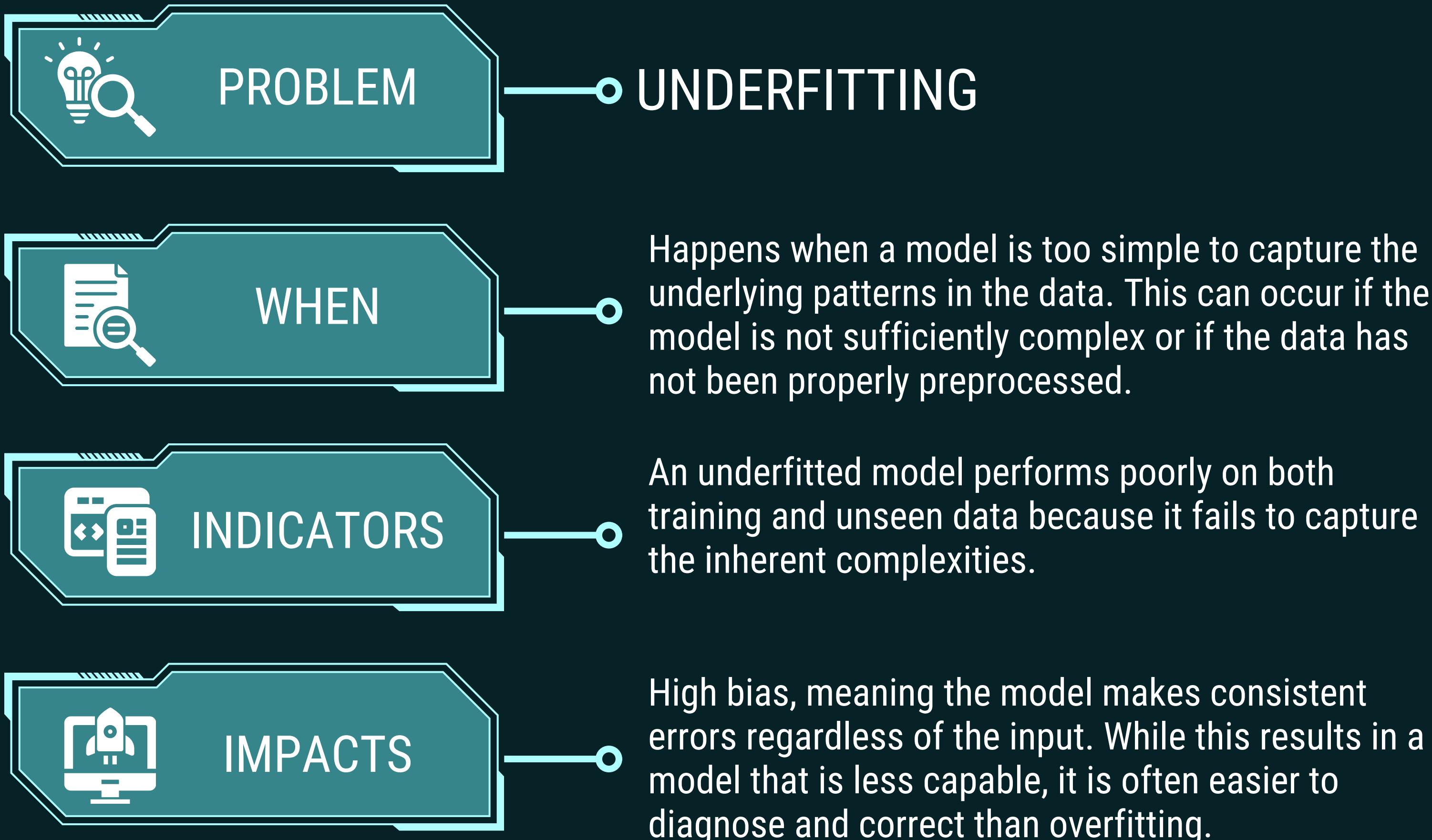




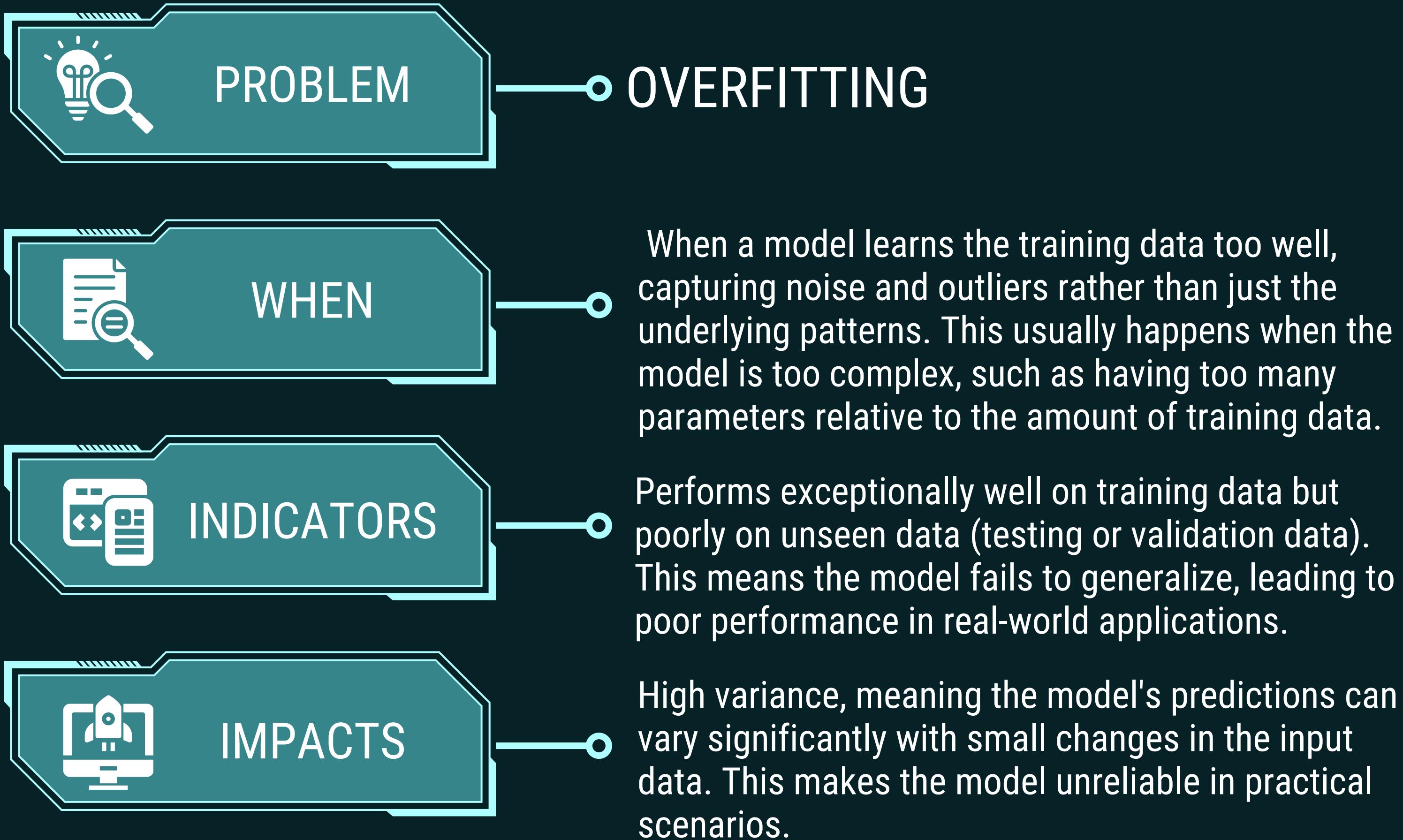
# OVERFITTING DETECTION & TROUBLESHOOTING



# OVERFIT VS. UNDERFIT

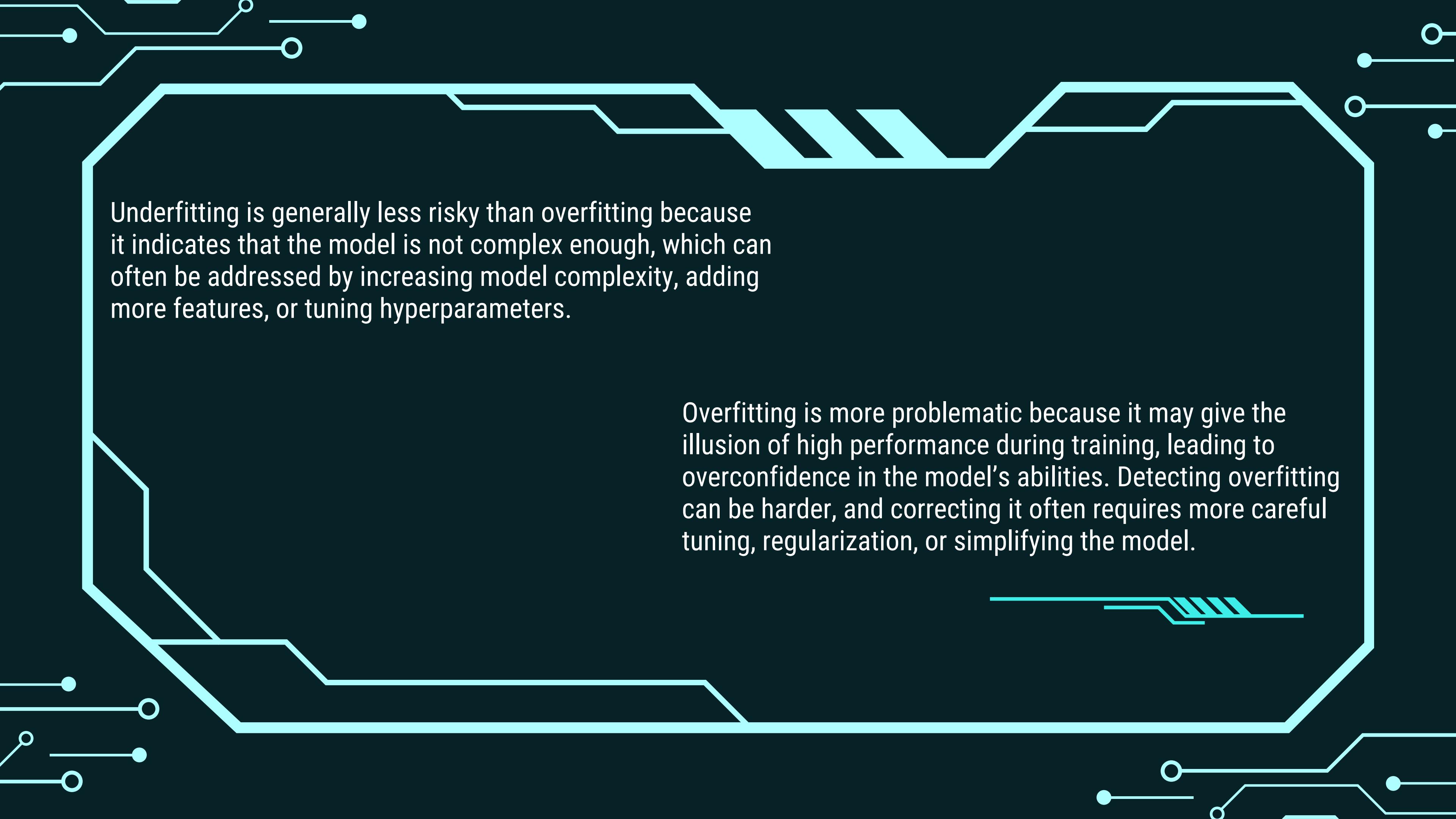


# OVERFIT VS. UNDERFIT





# WHICH IS LESS RISKY?



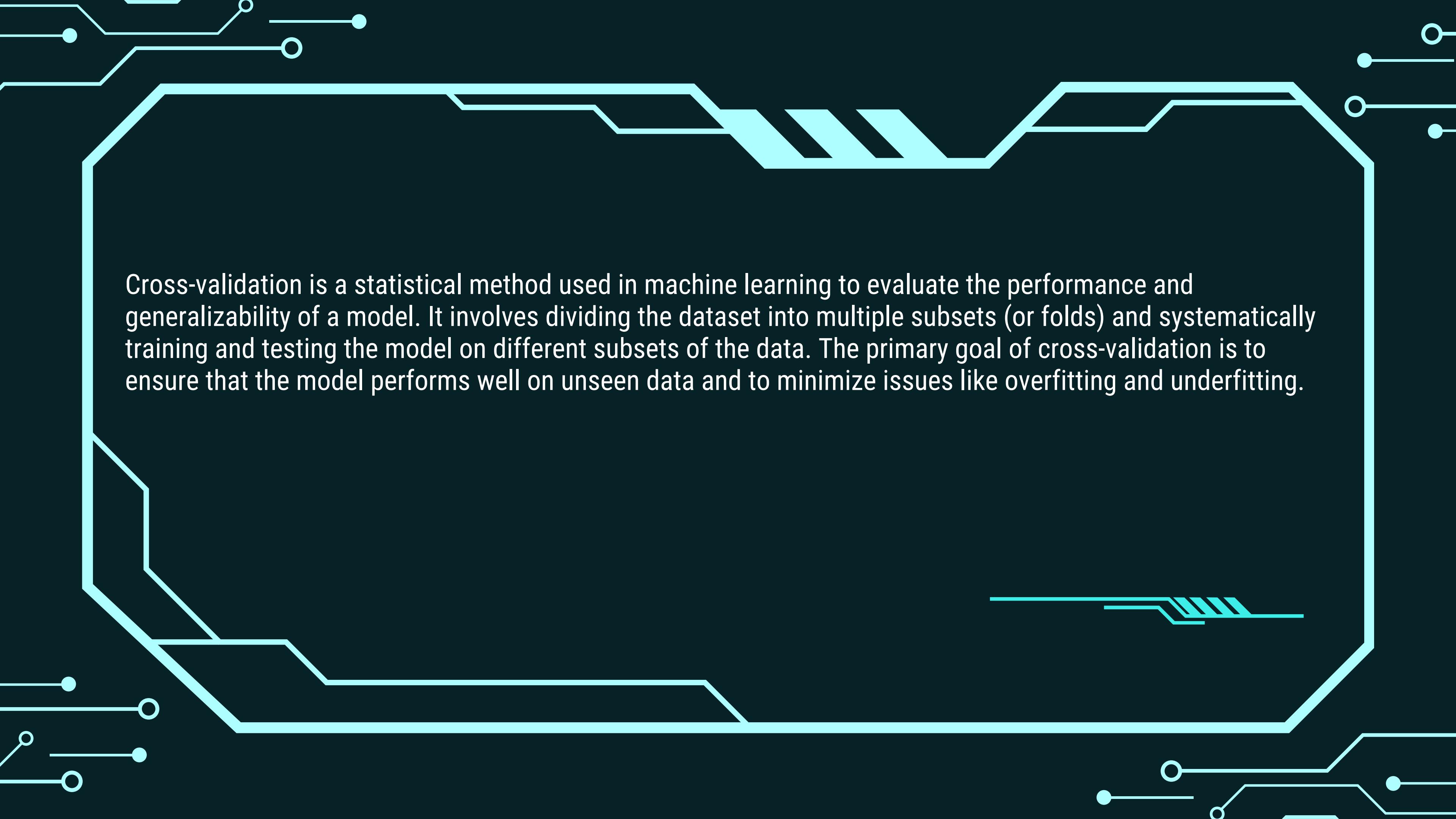
Underfitting is generally less risky than overfitting because it indicates that the model is not complex enough, which can often be addressed by increasing model complexity, adding more features, or tuning hyperparameters.

Overfitting is more problematic because it may give the illusion of high performance during training, leading to overconfidence in the model's abilities. Detecting overfitting can be harder, and correcting it often requires more careful tuning, regularization, or simplifying the model.



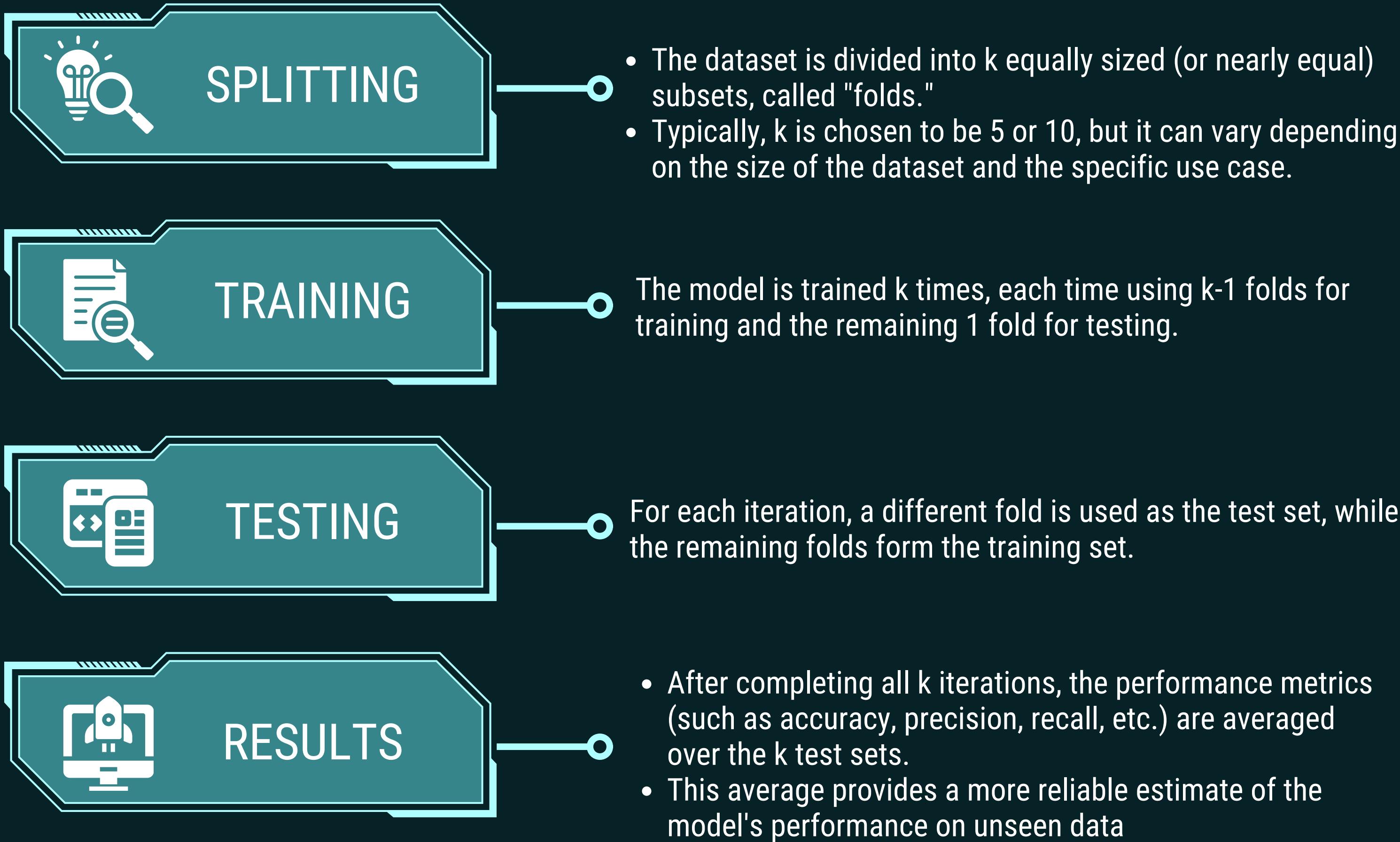
# TROUBLESHOOT CROSS-VALIDATION

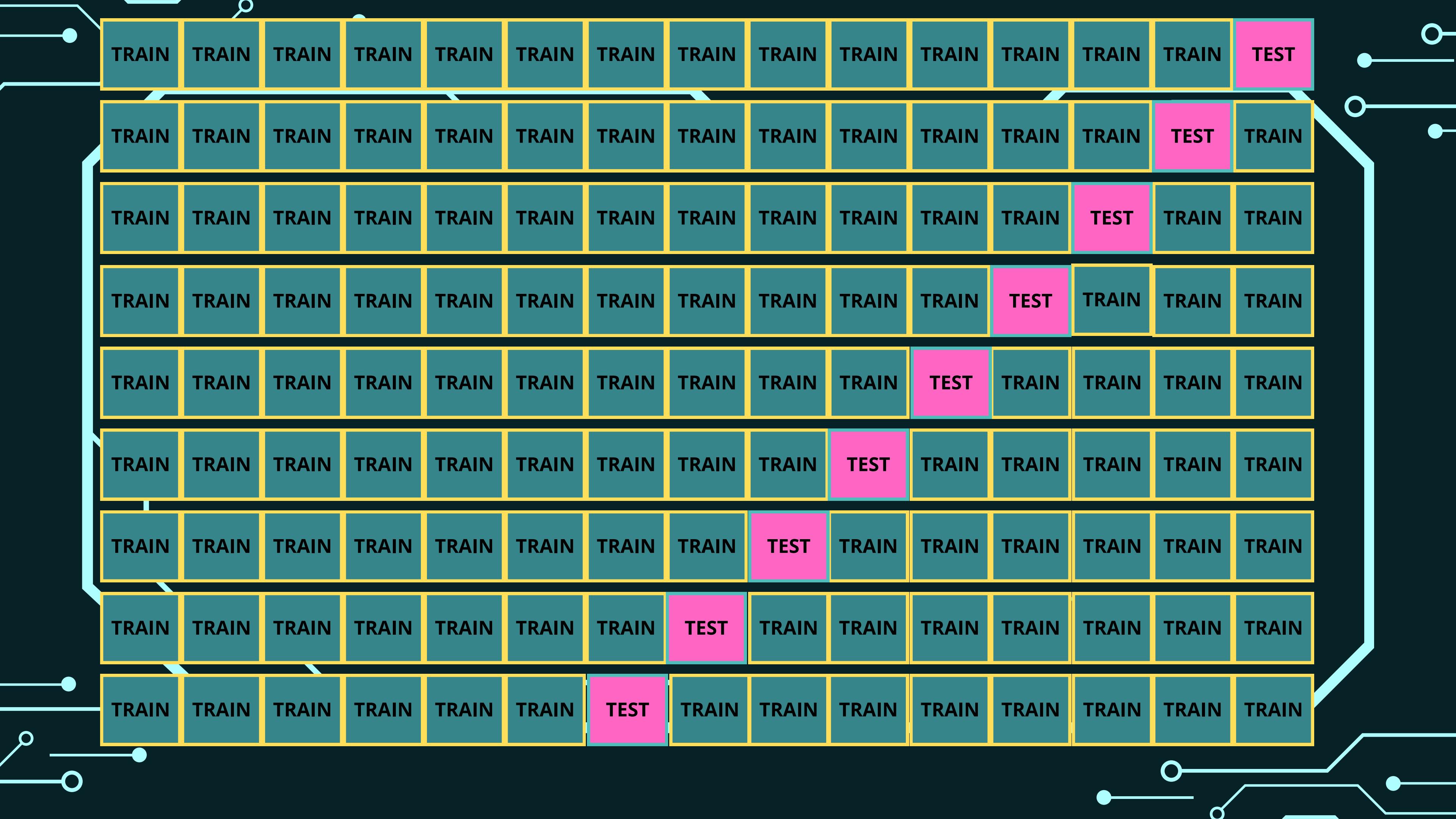




Cross-validation is a statistical method used in machine learning to evaluate the performance and generalizability of a model. It involves dividing the dataset into multiple subsets (or folds) and systematically training and testing the model on different subsets of the data. The primary goal of cross-validation is to ensure that the model performs well on unseen data and to minimize issues like overfitting and underfitting.

# OVERFIT VS. UNDERFIT





**ASSUMPTION:** Training, Test and Validation sets are independent of each other.

**Which of the two scenarios violate the assumption?**

- 1. Predicting the age of a person using facial features of people, in the data there are picture of identical twins, one was used in training, the other is used in testing.**
- 2. Predicting criminal activities of streets using house features, on the training set are houses at the right side of the road, on the testing set the houses on the left side of the road.**

# GENERALIZATION BOUNDARIES



## GENERALIZATION

The trained model should perform fairly well in new datasets.

## GENERALIZATION BOUNDARIES

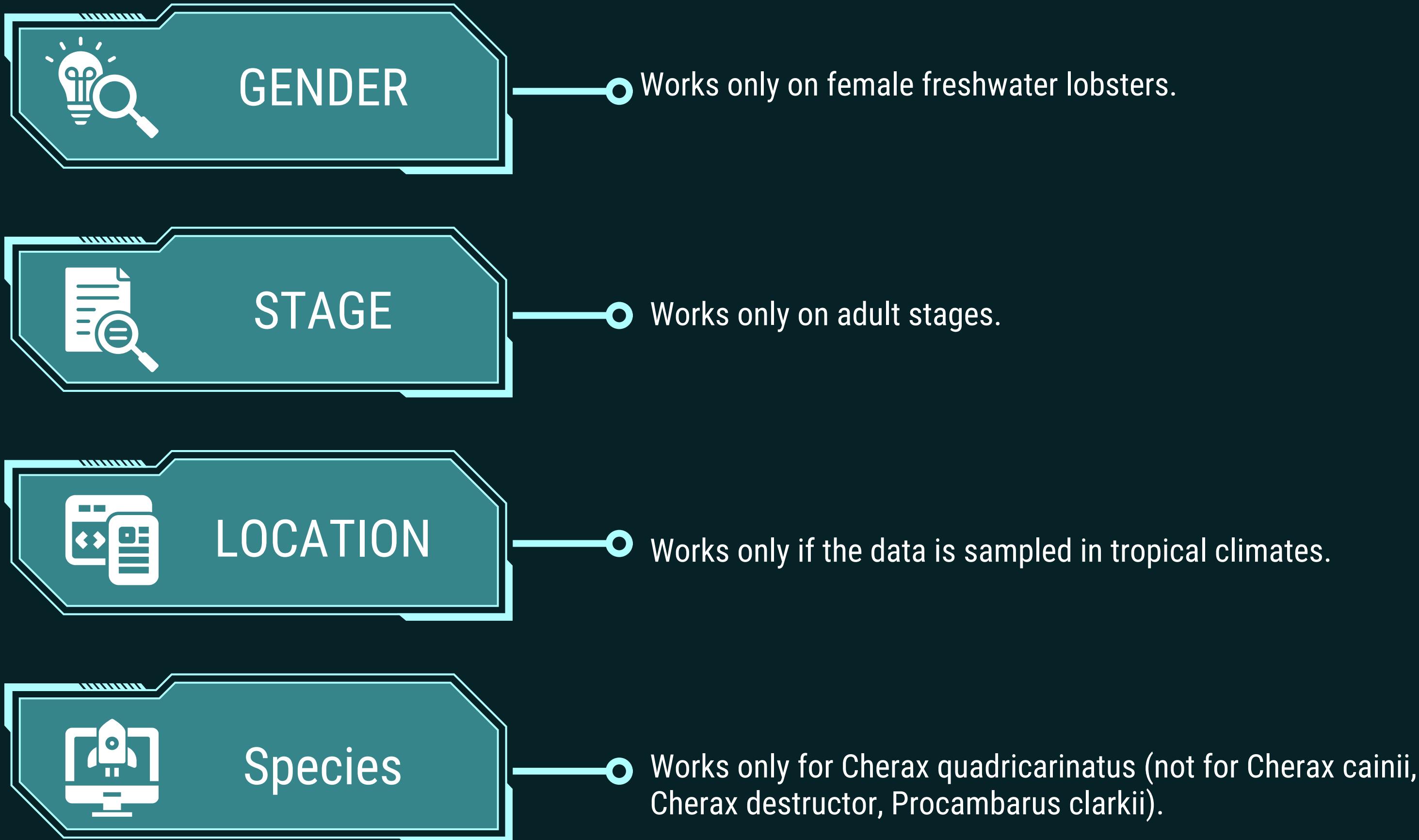
- Affected by your sampling methods.
- Affected by the innate characteristics of your data.

# Current Research

Aquaculture Biotechnology: Effects of Feed Crude Protein, Water Temperature, and Dissolved Oxygen on Fecundity of Female Freshwater Lobsters (*Charex quadricarinatus*).



# GENERALIZATION BOUNDARIES



# KEY TAKE-AWAYS



- Model is as good as the data it is trained on.
- Choose what kind of features you want to train your models with.
- Be transparent with the features you trained into your model.
- Your model will work good based on the generalization boundaries, but its performance will decrease if your data starts to drift from what it is trained on.

# THANK YOU!

A presentation is a method of sharing information, ideas, or arguments with an audience using spoken words and visual aids. It typically involves a speaker conveying content to inform, persuade, or entertain the listeners.

[www.linkedin.com/in/gerard-ompad](https://www.linkedin.com/in/gerard-ompad)

gerard.ompad@gmail.com

<https://github.com/gerard-ompad>