# On Deploying Machine Learners into Embedded Systems

**Tommaso Zoppi**, Andrea Ceccarelli, Andrea Bondavalli
RCL Group – University of Florence - Italy
e-mail: tommaso.zoppi@unifi.it

RCL
RESILIENT COMPUTING LAB

UNIVERSITÀ DEGLI STUDI FIRENZE
DIMAI
DIPARTIMENTO DI MATEMATICA E INFORMATICA "ULISSE DINI"

# Tabular Data

► Embedded and General-Purpose systems often share the need of analysing tabular data

– Features: system indicators (mainly networks)

– Label: normal behavior or specific type of attack

# What Anomalies are?

**Anomaly detection refers to the problem of finding patterns in data that do not conform to an expected behaviour[1]**



[1] Chandola, Varun, Arindam Banerjee, and Vipin Kumar. "Anomaly detection: A survey." ACM computing surveys (CSUR) 41.3 (2009): 15.

# Purpose of Anomaly Detectors

► Anomalies may have many root causes

– Security threats

– Misconfigurations

– Performance Issues

– Wrong/Slow interactions with other devices

– Benign alterations

► Regardless of their root cause, it is always beneficial to detect them.

# Embedding Anomaly Detectors

► Anomaly detectors usually rely on supervised/unsupervised ML algorithms

  – Which are usually resource and time-consuming

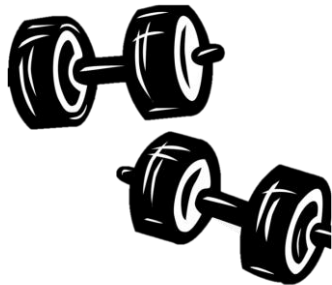  – Not a huge problem for systems that do not have hardware or real-time constraints

## BUT BUT BUT

► There always exists some kind of limitation to develop systems "in practise"

  – Thus, assuming "unlimited resources" is not doable

RCL
RESILIENT COMPUTING LAB

UNIVERSITÀ
DEGLI STUDI
FIRENZE
DIMAI
DIPARTIMENTO DI
MATEMATICA E INFORMATICA
"ULISSE DINI"

# Then what?

► As a result, the best intrusion/error/anomaly detector or failure predictor for a given system must be chosen according to constraints:
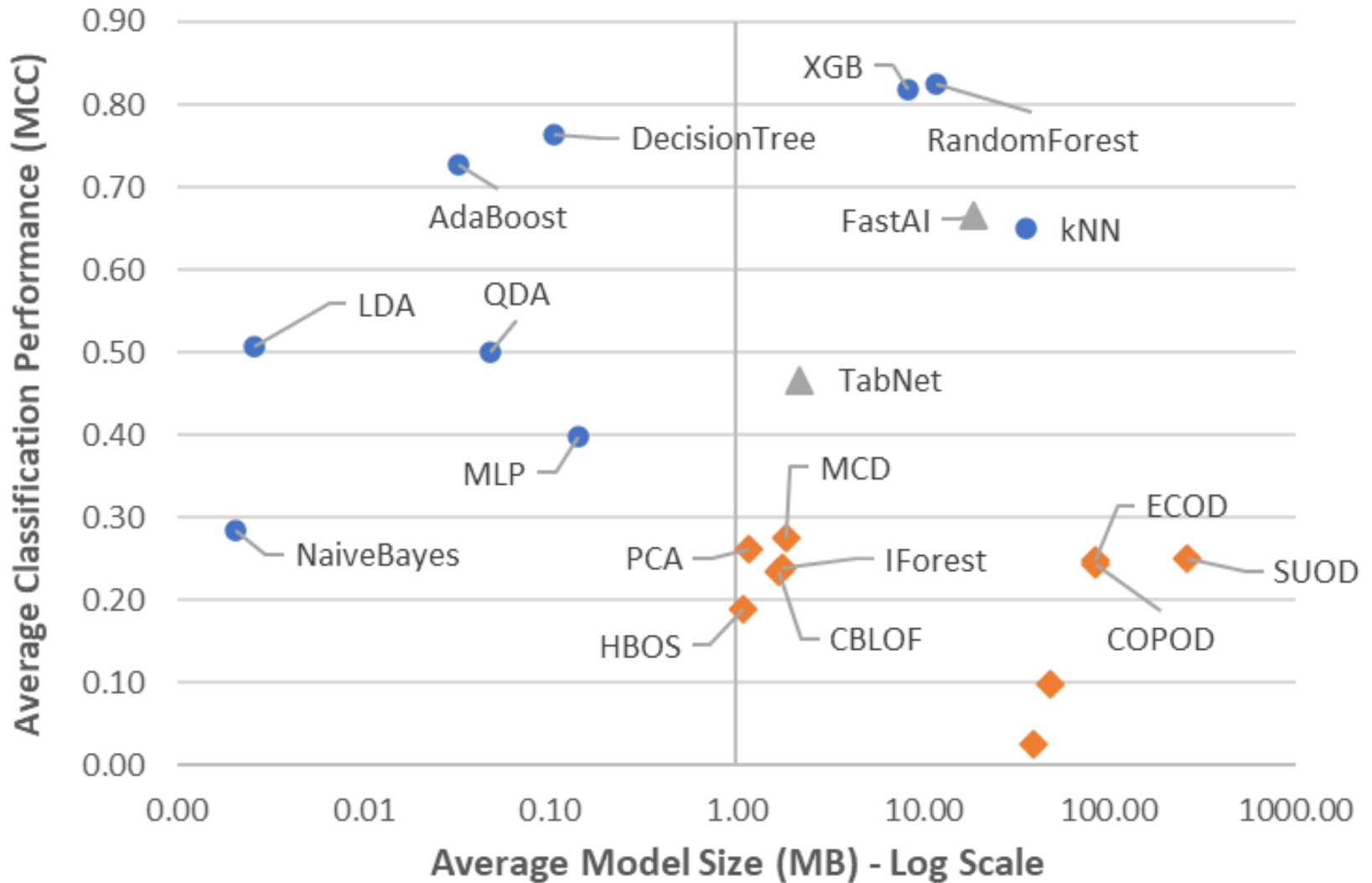
- Model size

- Model speed

- Detection performance
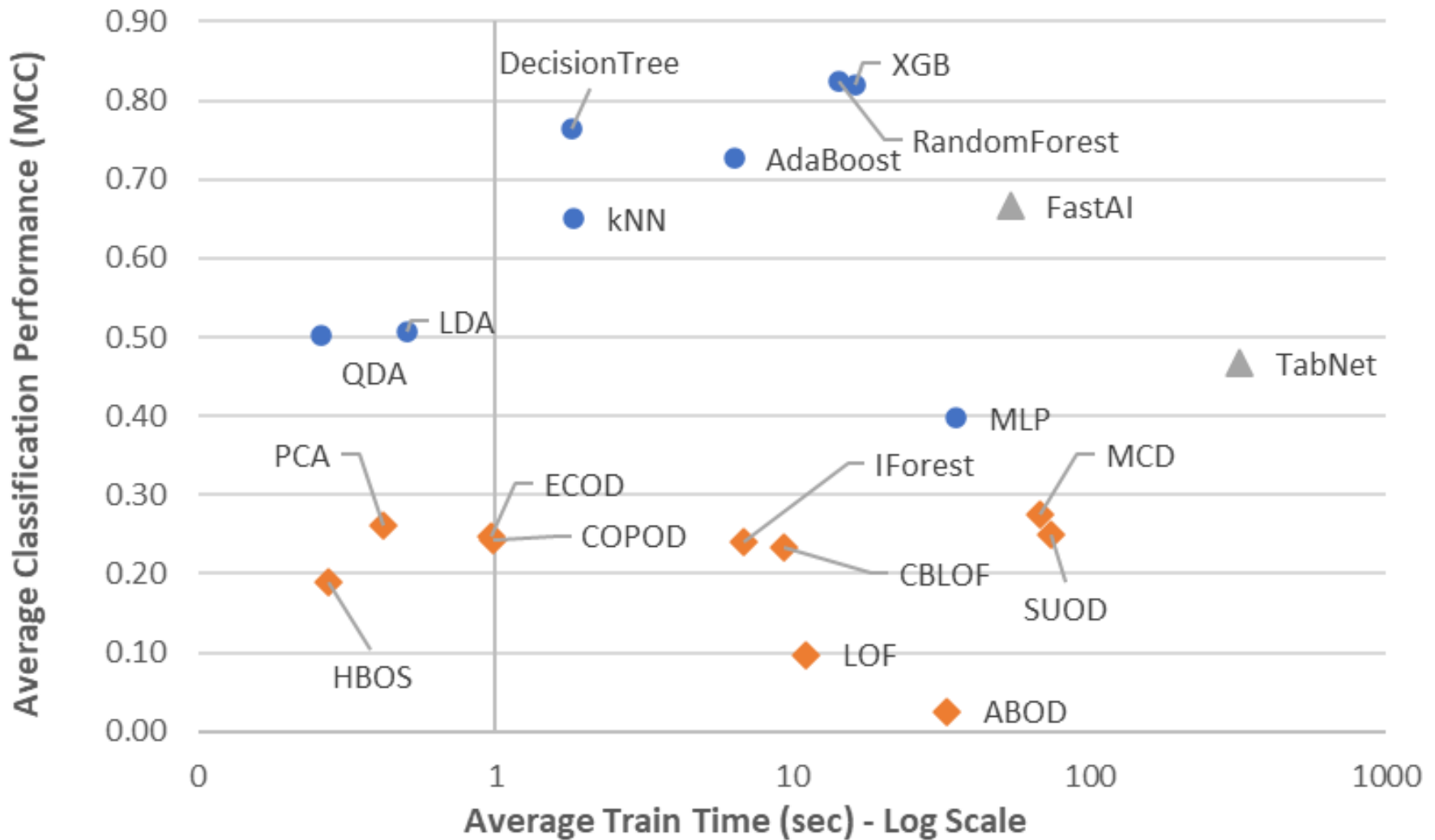
- Availability of labelled data for training

- ...

► That is why we took several SotA ML algorithms

- Supervised: DecisionTree, RandomForest, XGB, NaiveBayes, LDA, kNN, MLP, AdaBoost, QDA

- Unsupervised: COPOD, ABOD, HBOS, MCD, PCA, ECOD, LOF, CBLOF, Iforest, SUOD

- Deep learning: TabNet, FastAI

► And we exercised them on a total of 33 datasets regarding critical systems to derive their average performance metrics

RCL
RESILIENT COMPUTING LAB
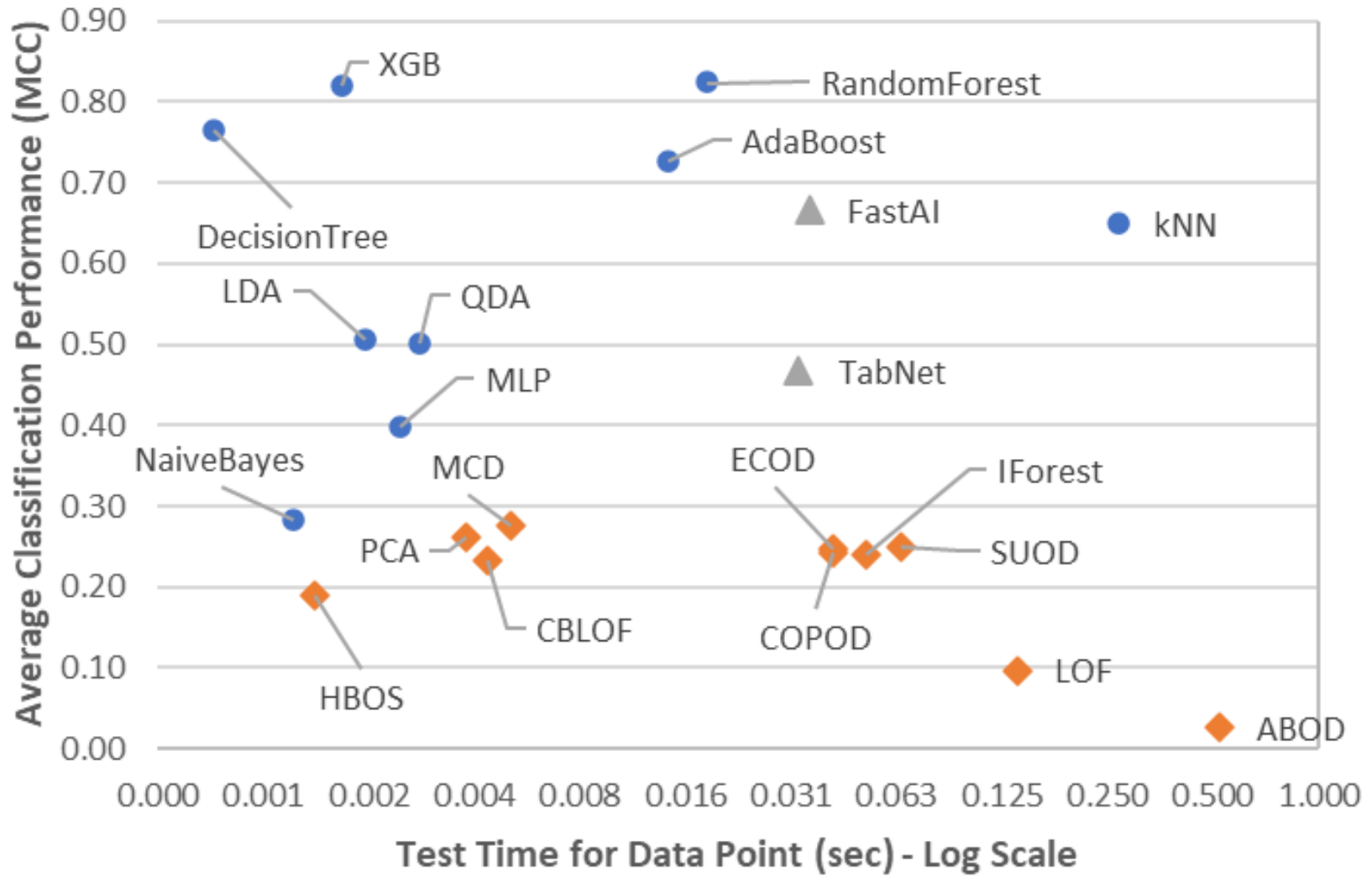
UNIVERSITÀ DEGLI STUDI FIRENZE
DIMAI
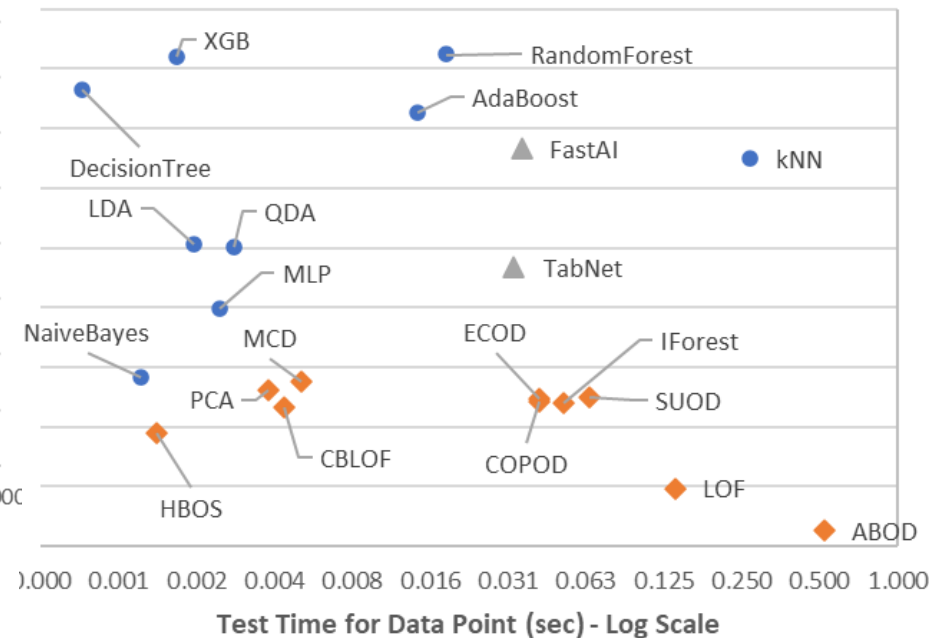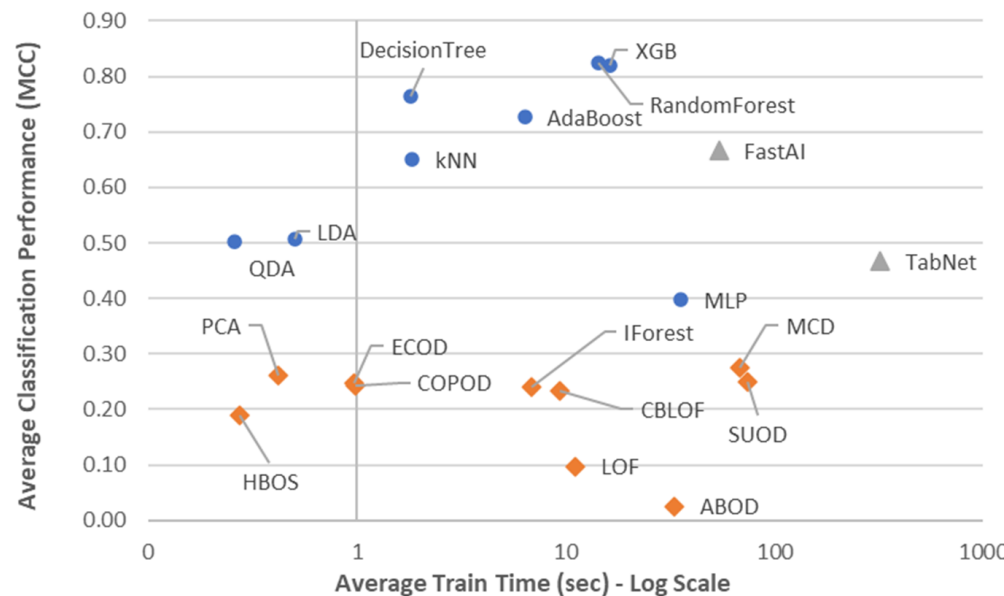DIPARTIMENTO DI MATEMATICA E INFORMATICA "ULISSE DINI"

# Model Size of Detectors

# Model Speed of Detectors (I)

# Model Speed of Detectors (II)

► **There are fast and slow algorithms**

- But also there are some that are fast during training and slow at runtime e.g., neighbour-based ones

- and vice versa

# Beware!

► Those numbers partially confirm the rather recent works stating that

**We should not think about deep learning as the panacea for any classification task!**

► For tabular data, tree-based classifiers are more interpretable, often faster and output fewer misclassifications

  – Good news for devices with limited resources!

  – See: Shwartz-Ziv, Ravid, and Amitai Armon. "**Tabular data: Deep learning is not all you need.**" Information Fusion 81 (2022): 84-90 (from AI ML group at Intel Israel)

# (Finally!) Wrapping Up…

► **This talk went through common constraints in deploying ML into embedded systems**

- There is no "silver bullet" algorithm to plug into a system for excellent detection capabilities and performance

- Detectors have to be crafted for specific systems depending on their constraints

- Availability of labels for training data is always scarce

- This calls for unsupervised detectors, which usually have poor detection capabilities

  - There are (research) works in the direction of making unsupervised ML more accurate

  - Get in touch with us if interested!

# Q&A Time