# Predicting Nanoparticle Toxicity Through Machine Learning-Based Models

Anne Christiono[1], Joseph Cave[2,3], Vittorio Cristini[2,3], Zhihui Wang[2,3], Prashant Dogra[2,3]

[1]William P. Clements High School, Sugar Land, TX; [2]Weill Cornell Graduate School of Medical Sciences, New York City, NY; [3]Mathematics in Medicine Program, HMRI, Houston, TX

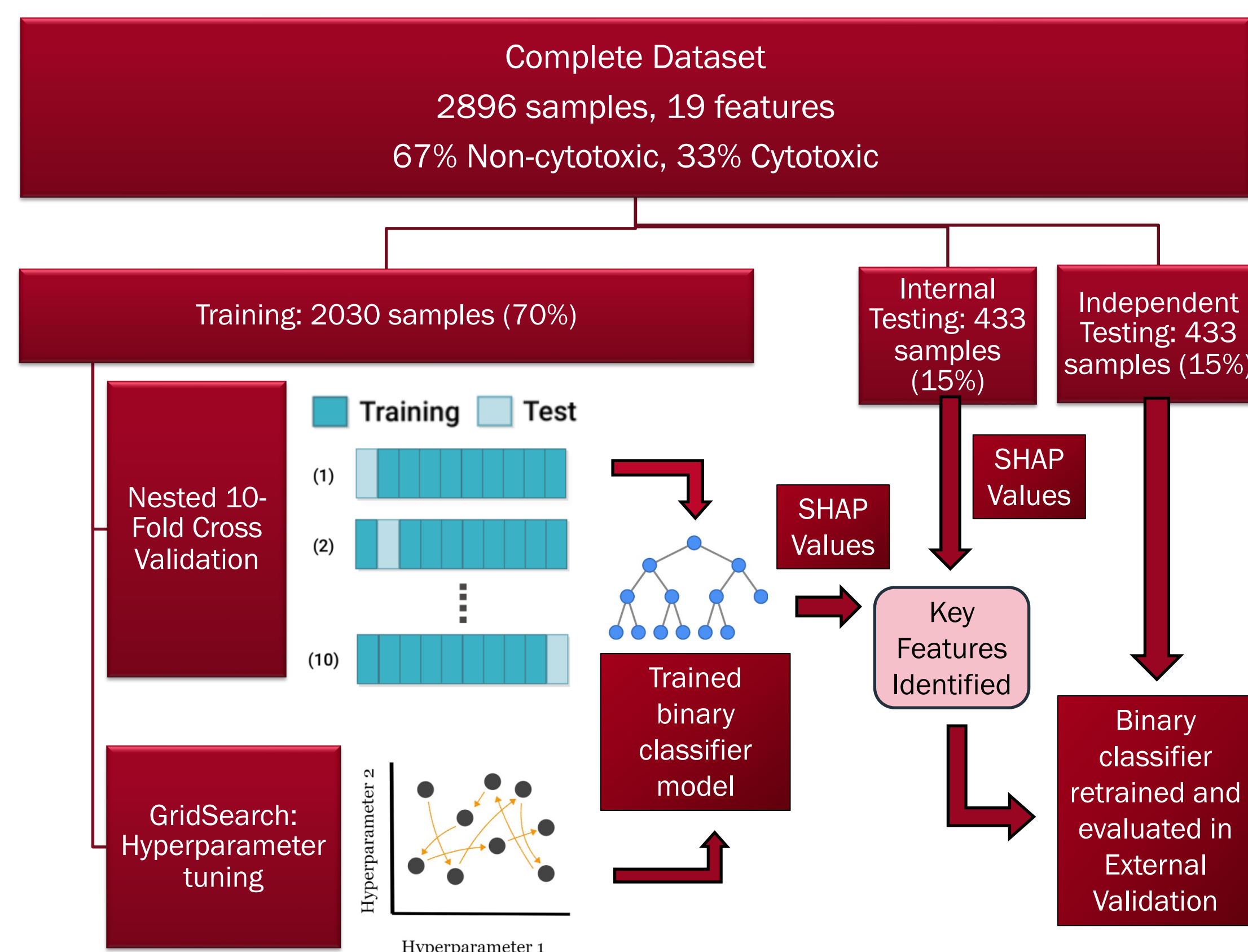HOUSTON Methodist® LEADING MEDICINE

## Introduction

- Engineered nanoparticles (NPs)—encompassing both inorganic NPs, such as metal oxides, and organic NPs, such as lipid-based NPs—are being leveraged in the design of effective and targeted drugs.
- Accurate prediction of NP cytotoxicity is critical to the design of safe NPs, but is also a challenge due to the complex toxicity mechanisms these NPs exhibit when interacting with biological environments.
- Machine learning can be leveraged in the development of reliable methods to determine NP cytotoxicity, ensuring the safe use of NPs in various medical applications and enabling efficient drug design.

## Purpose

We sought to develop a generalized machine learning model to accurately classify various types of NPs as cytotoxic or non-cytotoxic, and determine the most significant attributes which can explain NP toxicity, therefore facilitating the intentional engineering of safe NPs. We hypothesize that physicochemical NP-related properties would significantly impact cytotoxicity.

## Methods

- Through an extensive review of existing literature, we have chosen a comprehensive dataset of 2896 NPs, evaluated over 19 physicochemical and experimental attributes and encompassing 33 different NP types widely used in medical products.
- 70% of the dataset was reserved for model training, 15% for internal validation, and 15% for independent testing.
- We built 32 predictive models, with algorithms ranging between linear, nonlinear, and tree-based classifiers.
- To maximize model performance on independent data, prevent overfitting, and improve generalizability, each model was trained using nested 10-fold cross validation and GridSearch.
- We used the Shapley Additive exPlanation (SHAP) values evaluated over the best-performing models to identify key features significant in determining NP toxicity.
- With the exclusive use of the selected key features, we retrained and tested the performance of the best-performing models, assessing the capacity of these features alone to accurately predict nanoparticle cytotoxicity.

Complete Dataset
2896 samples, 19 features
67% Non-cytotoxic, 33% Cytotoxic

Training: 2030 samples (70%)
Internal Testing: 433 samples (15%)
Independent Testing: 433 samples (15%)

Nested 10-Fold Cross Validation

GridSearch: Hyperparameter tuning

Trained binary classifier model

SHAP Values

Key Features Identified

SHAP Values

Binary classifier retrained and evaluated in External Validation

## Results

**A)**

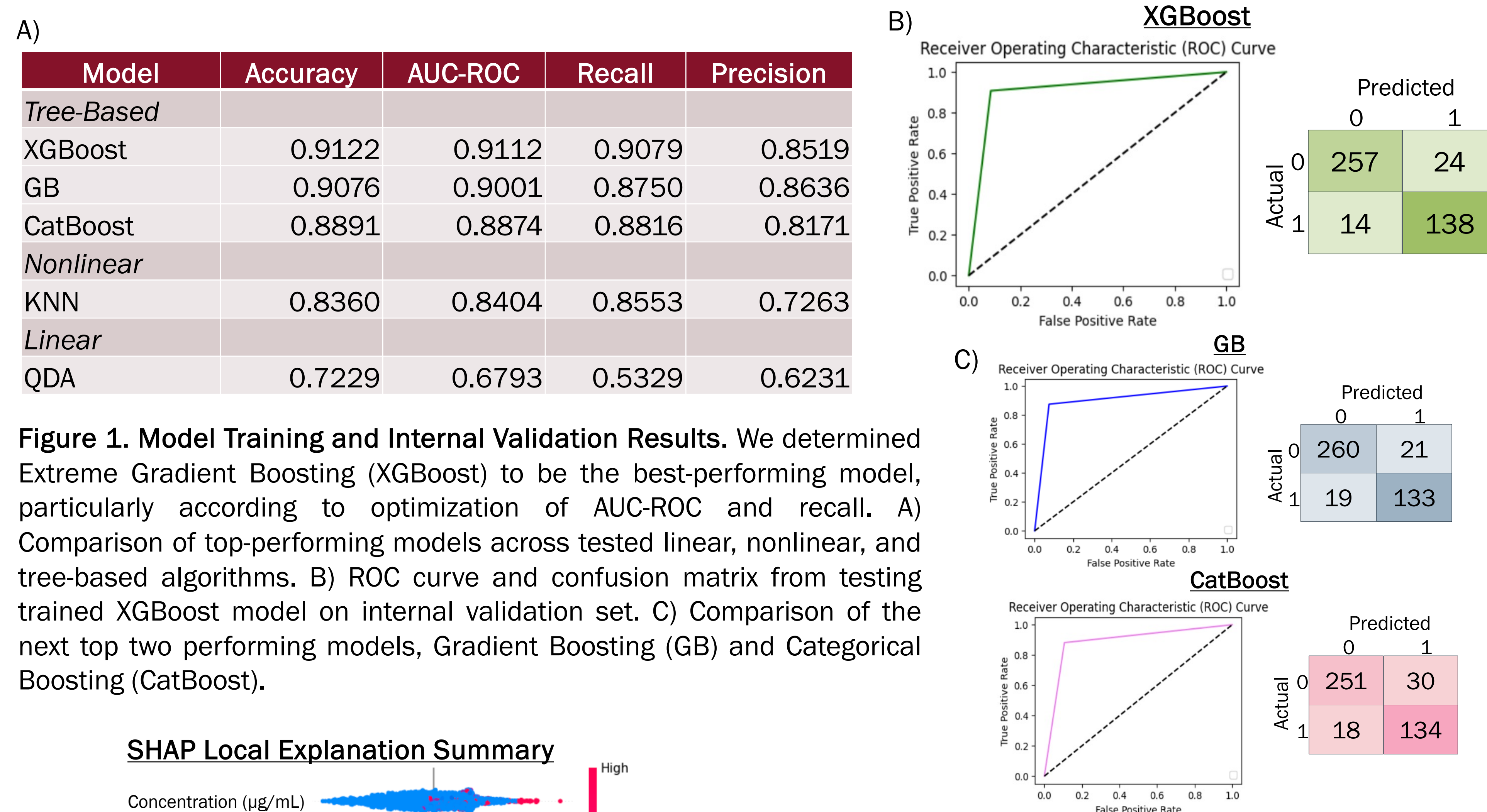| Model | Accuracy | AUC-ROC | Recall | Precision |
|---|---|---|---|---|
| *Tree-Based* | | | | |
| XGBoost | 0.9122 | 0.9112 | 0.9079 | 0.8519 |
| GB | 0.9076 | 0.9001 | 0.8750 | 0.8636 |
| CatBoost | 0.8891 | 0.8874 | 0.8816 | 0.8171 |
| *Nonlinear* | | | | |
| KNN | 0.8360 | 0.8404 | 0.8553 | 0.7263 |
| *Linear* | | | | |
| QDA | 0.7229 | 0.6793 | 0.5329 | 0.6231 |

**Figure 1. Model Training and Internal Validation Results.** We determined Extreme Gradient Boosting (XGBoost) to be the best-performing model, particularly according to optimization of AUC-ROC and recall. A) Comparison of top-performing models across tested linear, nonlinear, and tree-based algorithms. B) ROC curve and confusion matrix from testing trained XGBoost model on internal validation set. C) Comparison of the next top two performing models, Gradient Boosting (GB) and Categorical Boosting (CatBoost).
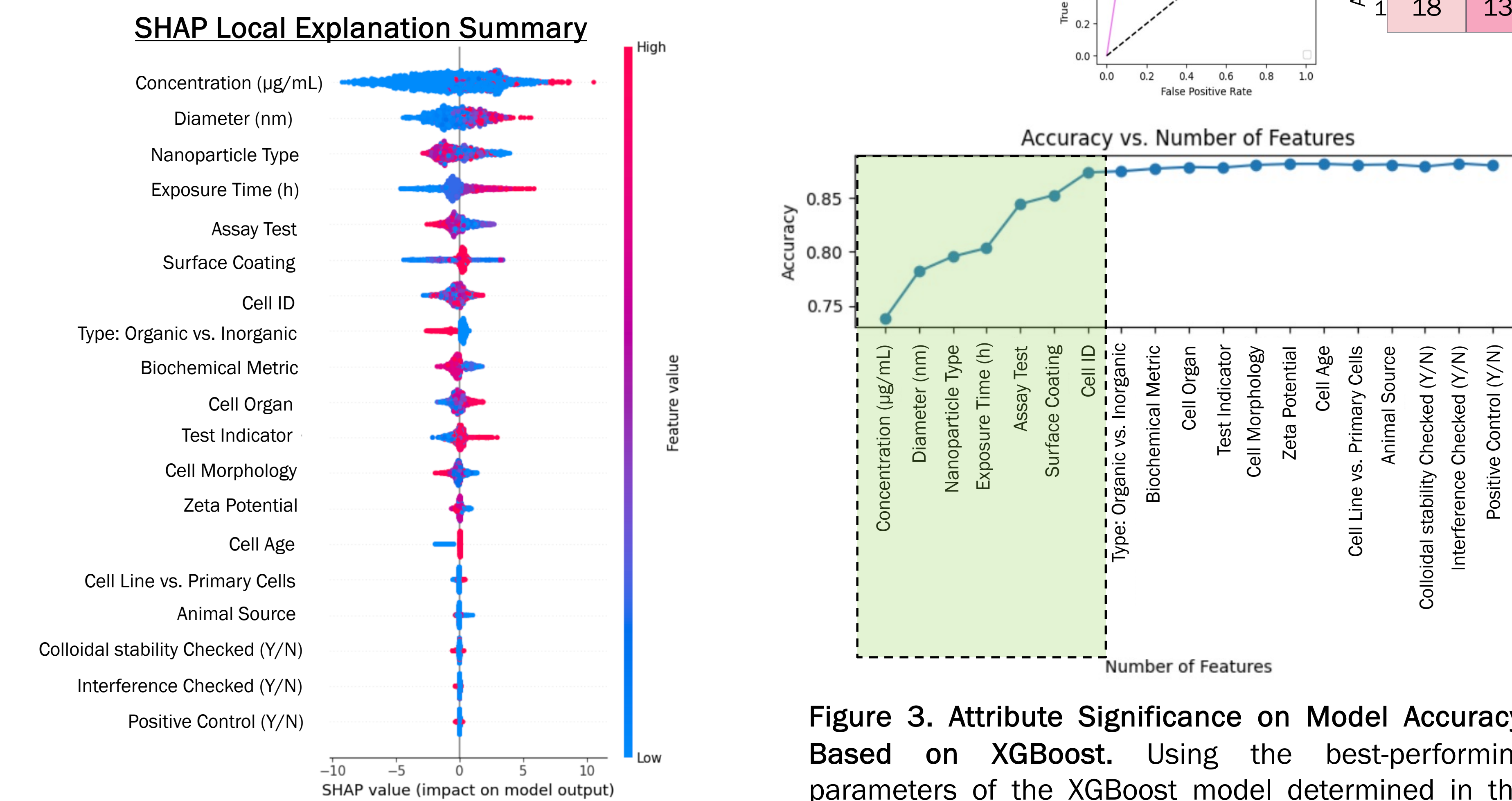
**B)**
XGBoost

**C)**
GB
CatBoost

**SHAP Local Explanation Summary**

Concentration (µg/mL)
Diameter (nm)
Nanoparticle Type
Exposure Time (h)
Assay Test
Surface Coating
Cell ID
Type: Organic vs. Inorganic
Biochemical Metric
Cell Organ
Test Indicator
Cell Morphology
Zeta Potential
Cell Age
Cell Line vs. Primary Cells
Animal Source
Colloidal stability Checked (Y/N)
Interference Checked (Y/N)
Positive Control (Y/N)

**Figure 2. SHAP Summary Plot.** Features are ranked based on their importance, evaluated by the SHAP metric. Each dot corresponds to an individual NP sample.
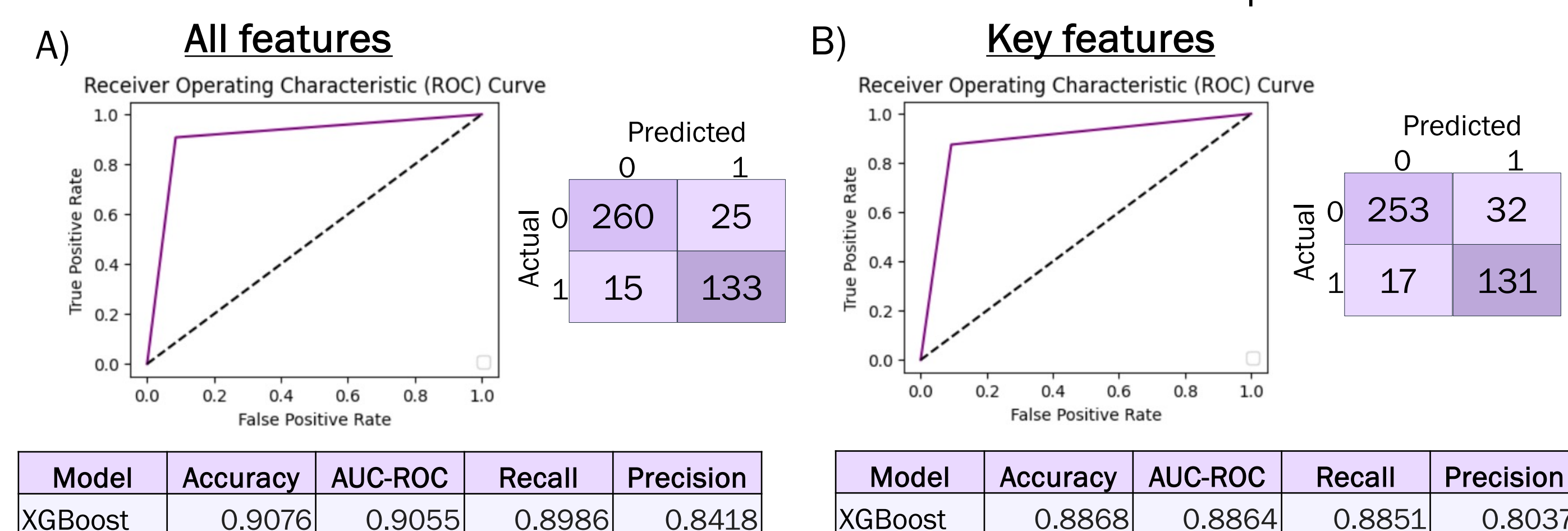
Accuracy vs. Number of Features

**Figure 3. Attribute Significance on Model Accuracy, Based on XGBoost.** Using the best-performing parameters of the XGBoost model determined in the training stage, features were incrementally incorporated one by one into the model. After each iteration, the model's accuracy was reevaluated to observe the impact of each added feature.

**A) All features**

| Model | Accuracy | AUC-ROC | Recall | Precision |
|---|---|---|---|---|
| XGBoost | 0.9076 | 0.9055 | 0.8986 | 0.8418 |

**B) Key features**

| Model | Accuracy | AUC-ROC | Recall | Precision |
|---|---|---|---|---|
| XGBoost | 0.8868 | 0.8864 | 0.8851 | 0.8037 |

**Figure 4. Retrained Model Tested on Independent Data .** A) XGBoost trained on all 19 features and tested on independent data (433 samples). B) XGBoost retrained on 7 key features, tested on same independent data.

## Results/Implications

- Top-performing models are able to classify NPs based on cytotoxicity at around 90% accuracy, with considerable emphasis on accurately identifying cytotoxic NPs (Fig 1).
- Developed models demonstrated consensus that NP cytotoxicity is primarily associated with parameters such as *NP concentration, size, NP type, exposure time, assay test type, surface coating, and specific cell type* (Fig 2).
- NP zeta potential is ranked as a feature with low importance, potentially attributed to its categorization as positive, negative, or neutral due to limitations of the dataset. To extend this work, we will recharacterize zeta potential as a continuous variables, and reassess its influence on cytotoxicity.
- Selected features encompass physicochemical and experimental properties, suggesting the potential for deliberate engineering and modification of NPs to mitigate toxicity risks.
- When comparing the model's accuracy considering all parameters to that achieved by solely considering the 7 most significant features, a similar level of performance is observed. The addition of the remaining features contribute to only marginal improvements in the model's performance in predicting NP cytotoxicity (Fig 3).
- Model effectively predicts NP cytotoxicity solely using the 7 most significant features with comparable accuracy to when all 19 features are considered, thereby illustrating the substantial influence of these selected features (Fig 4).
- This study also showcases the model's capacity to make accurate predictions across diverse types of NPs, emphasizing its robustness and versatility in NP characterization and toxicity assessment.

## Future Actions

- To enhance the model's performance and generalizability, we will augment the dataset by collecting additional data from relevant literature sources.
- We will recast certain parameters, including recharacterizing zeta potential as a numerical variable and normalizing size using logarithm transformation.
- We will conduct an in-depth analysis of the significant features identified by the models, aiming to discern precise nanoparticle modifications that can deliberately reduce cytotoxicity risk.
- To establish a consensus on the key attributes identified by the SHAP metric, we will train additional models, particularly those with embedded feature selection algorithms, and utilize different scoring metrics to further validate the significance of the selected features.
- We will expand the pipeline to address NP biodistribution behavior, allowing for a more comprehensive understanding of the NP-related factors influencing the biodistribution of NPs specifically targeted to tumor sites.

## Acknowledgments

## References

1. Labouta, H et al. Meta-Analysis of Nanoparticle Cytotoxicity via Data-Mining the Literature. *ACS Nano*, 2019. doi: 10.1021/acsnano.8b07562.
2. Martin, R et al. Evidence-Based Prediction of Cellular Toxicity for Amorphous Silica Nanoparticles. *ACS Nano 17*, no. 11 (2023): 9987–99. doi: 10.1021/acsnano.2c11968.