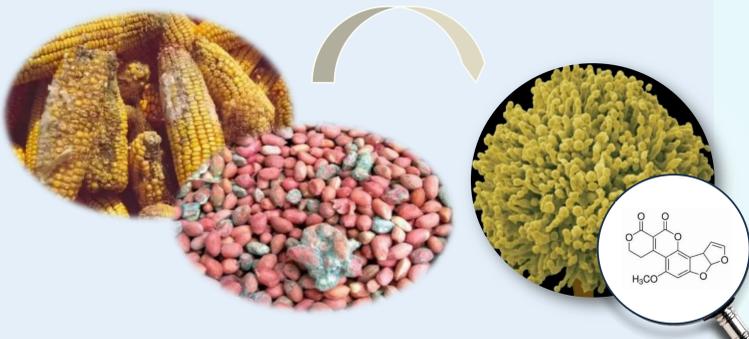


## I. Introduction

*Aspergillus flavus* is a major cause of post-harvest losses in maize and peanuts due to production of aflatoxin B<sub>1</sub> (AFB<sub>1</sub>) contamination. Current detection methods are often slow, invasive, and unsuitable for real-time monitoring, highlighting the urgent need for sensitive and scalable early detection strategies.

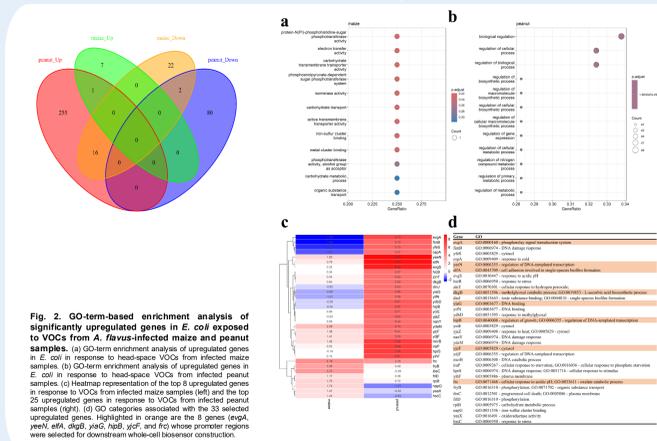
## II. Objectives

This study aims to develop a non-invasive, transcriptome-guided whole-cell biosensor array integrated with machine learning models for early and quantitative prediction of *A. flavus* infection stages and AFB<sub>1</sub> levels in maize and peanut kernels.



## IV. Results and Expectation

We identified eight infection-induced promoters in *E. coli* by analyzing its transcriptomic response to HVOCs in infected maize and peanuts.



Feature importance analysis revealed that early host responses, including transcriptional regulation and biofilm formation, served as key predictive features, thereby providing mechanistic interpretability not attainable with conventional optical or chemical assays.

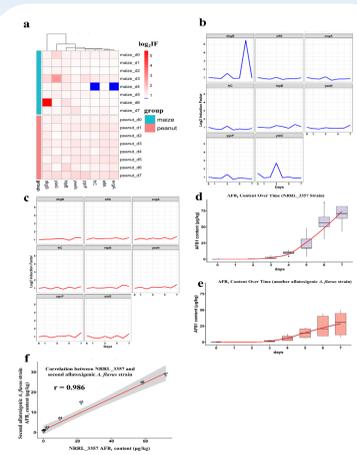


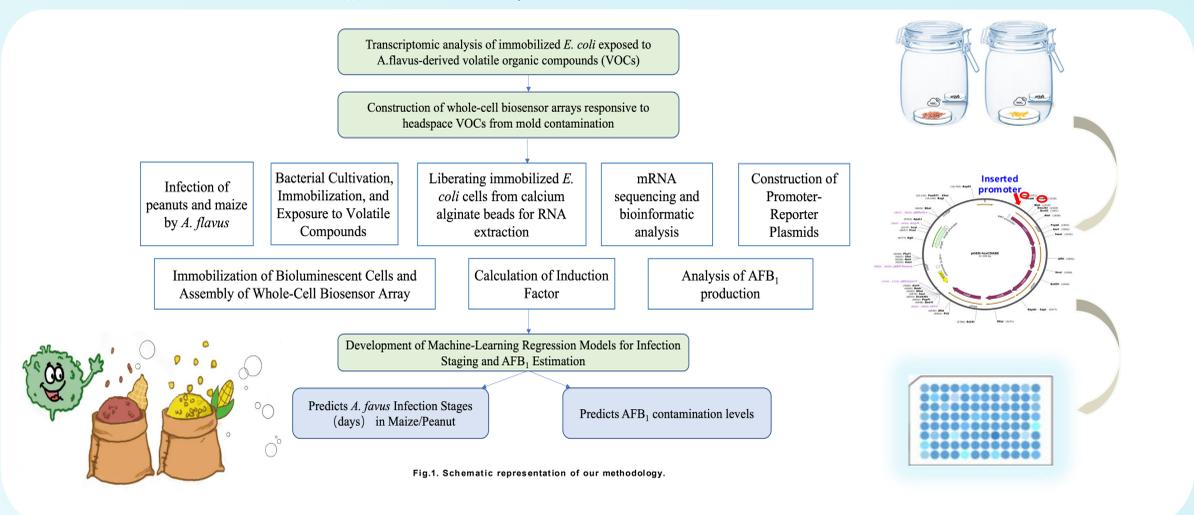
Fig. 3. Temporal dynamics and clustering of whole-cell bioreporter responses to *A. flavus*-induced VOCs in maize and peanut. (a) Hierarchical clustering heatmap of Log<sub>2</sub>-transformed induction factor (Log<sub>2</sub>IF) values across 8 promoter-driven biotransformers exposed to VOCs emitted from *A. flavus*-infected kernels, illustrating distinct temporal activation profiles. (b) Time-course plots of Log<sub>2</sub>IF values in maize samples (DMS) measured at 0, 1, 2, 3, 4, 5, 6, and 7 days post-inoculation (DPI). Each panel corresponds to a specific promoter. (c) Corresponding Log<sub>2</sub>IF trajectories in peanut samples over the same infection timeline. (d) AFB<sub>1</sub> accumulation dynamics in maize kernels inoculated with *A. flavus* strain NRRL\_3357 (c) and a second aflatoxigenic strain (e), quantified across identical DPI stages. (f) Pearson correlation coefficient of AFB<sub>1</sub> levels between the two strains across all time points, indicating concordance in temporal AFB<sub>1</sub> accumulation patterns.

## V. Conclusion

This integrated biosensor-machine learning platform provides a robust, non-invasive approach for real-time monitoring of *A. flavus* contamination, offering significant potential for scalable food safety applications across diverse agroecosystems.

## III. Methods

A transcriptome-guided whole-cell biosensor array was developed by integrating eight infection-induced promoters, identified from *E. coli* transcriptomic responses to the headspace volatile organic compounds (HVOCs), into calcium alginate-immobilized bioreporters coupled with machine learning regression models, as summarized in Figure 1. Time-resolved bioluminescence signals were used to train ensemble regressors, including XGBoost, CatBoost, and RandomForest, for quantitative prediction of infection stages and AFB<sub>1</sub> levels.



For maize: XGBoost consistently achieved superior performance in internal validation, with R<sup>2</sup> values of 0.94 and 0.98 for maize infection staging and AFB<sub>1</sub> quantification, respectively.

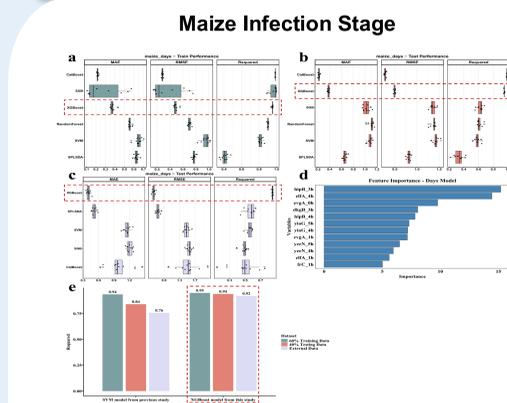


Fig. 4. Comparative performance of regression models for predicting *A. flavus* infection stages in maize kernels. (a-c) Boxplots summarizing the predictive accuracy of six regression algorithms—RF, SGB, CatBoost, XGBoost, SVM, and SPLSDA—evaluated using mean absolute error (MAE), root mean square error (RMSE), and coefficient of determination (R<sup>2</sup>). Models were trained and tested on: (a) 60% of maize kernel samples inoculated with *A. flavus* strain NRRL\_3357 (training set), (b) 40% of the same dataset (testing set), and (c) an independent external maize kernel dataset inoculated with a distinct aflatoxigenic *A. flavus* strain. (d) Feature importance plot showing the top 20 variables contributing to infection stage prediction, as identified by the best-performing XGBoost model, with importance scaled from 0 to 100. (e) Parallel comparison of R<sup>2</sup> values between an SVM model which utilized pre-validated 14 promoters for whole-cell biosensor array construction, and the current XGBoost model informed by transcriptome-guided selection of 8 promoters in this study, revealing enhanced predictive accuracy in the latter across testing (40%) and external datasets. The red dashed rectangles in panels (a-c) highlight the most accurate and robust models across metrics and datasets.

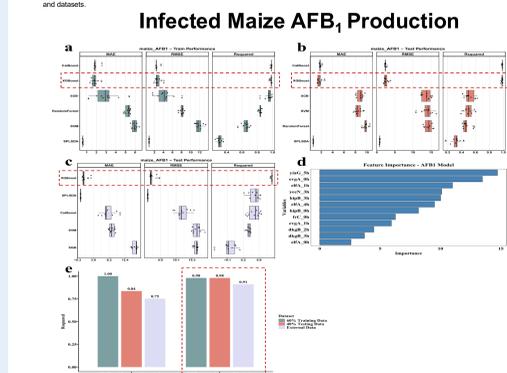


Fig. 5. Predictive performance benchmarking of regression models for estimating AFB<sub>1</sub> production by *A. flavus* in infected maize AFB<sub>1</sub> production. (a-c) Boxplots illustrating the predictive accuracy of six machine learning regression algorithms—RF, SGB, CatBoost, XGBoost, SVM, and SPLSDA—evaluated using MAE, RMSE, and R<sup>2</sup>. Models were trained and validated on: (a) 60% of peanut kernel samples inoculated with *A. flavus* strain NRRL\_3357 (training set), (b) 40% of the same dataset (testing set), and (c) an independent external maize kernel dataset inoculated with a genetically distinct aflatoxigenic *A. flavus* strain. (d) Feature importance plot from the top-performing CatBoost model, highlighting the 20 most influential variables contributing to model predictions, with normalized importance scores (0-100%). (e) Parallel comparison.

## VI. Conclusion

The model also maintained strong generalization in external validation using independent *A. flavus* isolates, achieving R<sup>2</sup> values of 0.92 and 0.91.

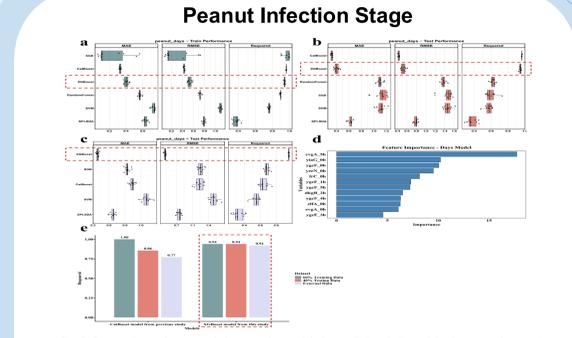
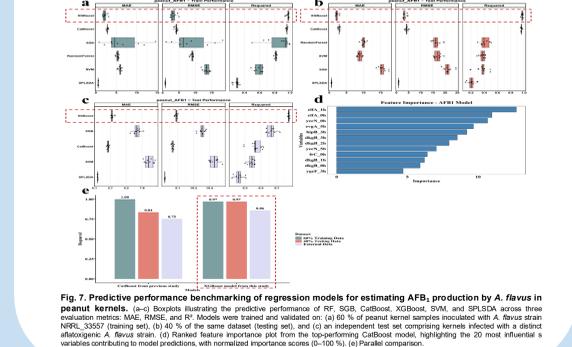


Fig. 6. Comparative performance of six regression models for predicting *A. flavus* infection stages in peanut kernels. (a-c) Boxplots illustrating the predictive performance of RF, SGB, CatBoost, XGBoost, SVM, and SPLSDA across three evaluation metrics: MAE, RMSE, and R<sup>2</sup>. Models were trained and validated on: (a) 60% of peanut kernel samples inoculated with *A. flavus* strain NRRL\_3357 (training set), (b) 40% of the same dataset (testing set), and (c) an independent test set comprising kernels infected with a distinct aflatoxigenic *A. flavus* strain. (d) Ranked feature importance plot from the top-performing CatBoost model, highlighting the 20 most influential variables contributing to model predictions, with normalized importance scores (0-100%). (e) Parallel comparison.



For peanuts: XGBoost achieved high accuracy in internal validation (R<sup>2</sup> = 0.94 and 0.97) and maintained robust performance in external validation (R<sup>2</sup> = 0.92 and 0.86).

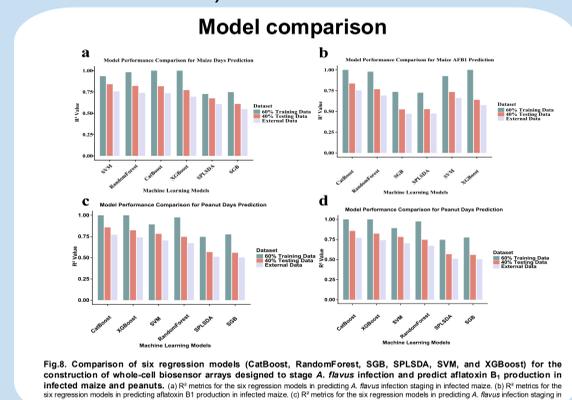


Fig. 8. Comparison of six regression models (CatBoost, RandomForest, SGB, SPLSDA, SVM, and XGBoost) for the construction of whole-cell biosensor arrays designed to stage *A. flavus* infection and predict aflatoxin B<sub>1</sub> production in infected maize and peanuts. (a) RF metrics for the six regression models in predicting *A. flavus* infection staging in infected maize. (b) RF metrics for the six regression models in predicting *A. flavus* infection staging in infected peanuts. (c) RF metrics for the six regression models in predicting aflatoxin B<sub>1</sub> production in infected maize. (d) RF metrics for the six regression models in predicting aflatoxin B<sub>1</sub> production in infected peanuts.