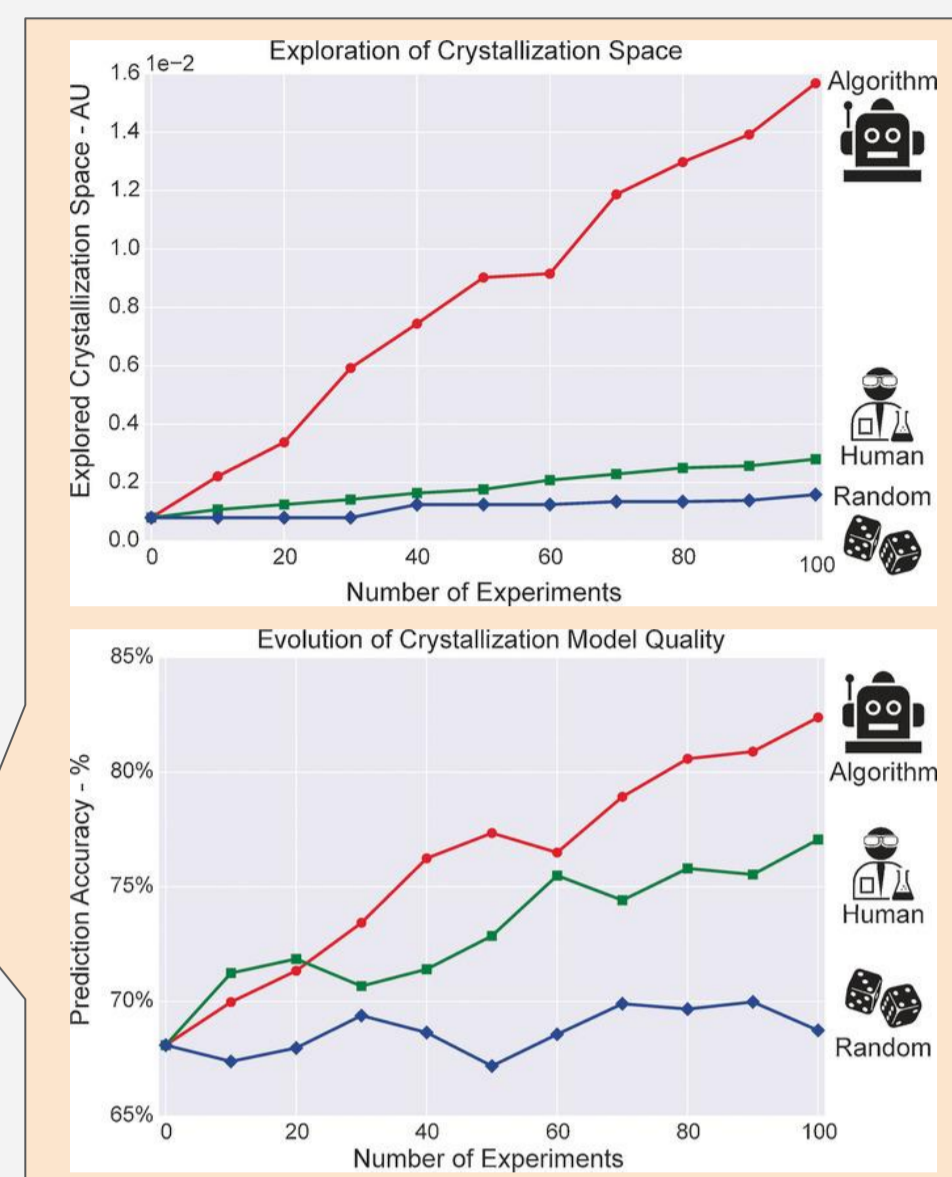
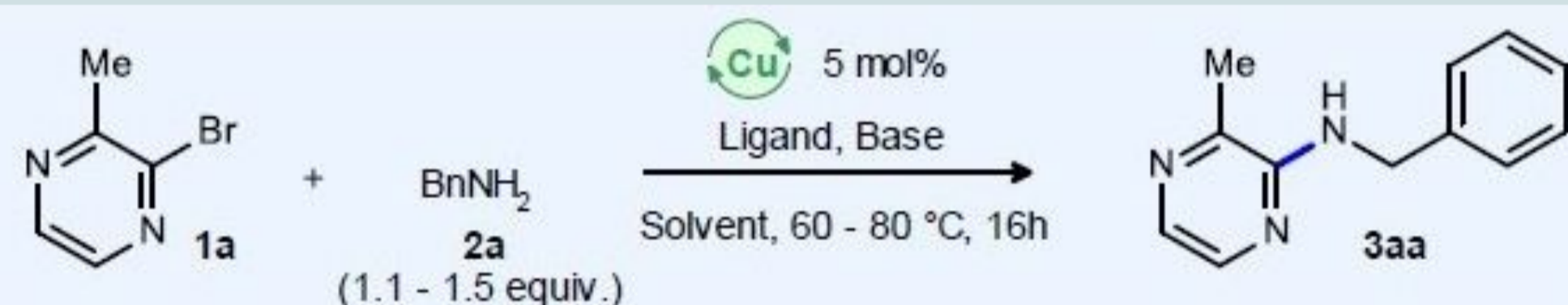


In recent years, **artificial intelligence (AI)** has been trialed as a tool to accelerate **chemical synthesis** as it can learn from existing experiments to propose new optimized reactions. It has been shown that certain robots are already capable of being as efficient as and even **outperforming** human chemists in select applications, but **how**? In this project, we will train our own **machine learning (ML)** model with a given data set and explain its predictions to encourage an increase in the usage of AI in chemical synthesis.

Braconi and Godineau studied how a **Bayesian Optimizer (BO)** could be used to discover high yielding reactions within a defined space. The model discovered reactions with yields up to **87.2%** in **264** attempts! Is it possible to find equal or better combinations with another method by using the data they provide?



### Choice of model ?

Several ML models are still being tested and developed for chemistry. A model that is **too simple** (i.e. linear regressor) will be **inefficient**, but a model that is **too complex** (i.e. Deep Neural Network) requires a lot of resources and is **hard to understand**...

**We want a model that is efficient but also interpretable. Can we guarantee both ?**



Using **XAI** and tree based ensemble methods, yes !



### What is XAI ?

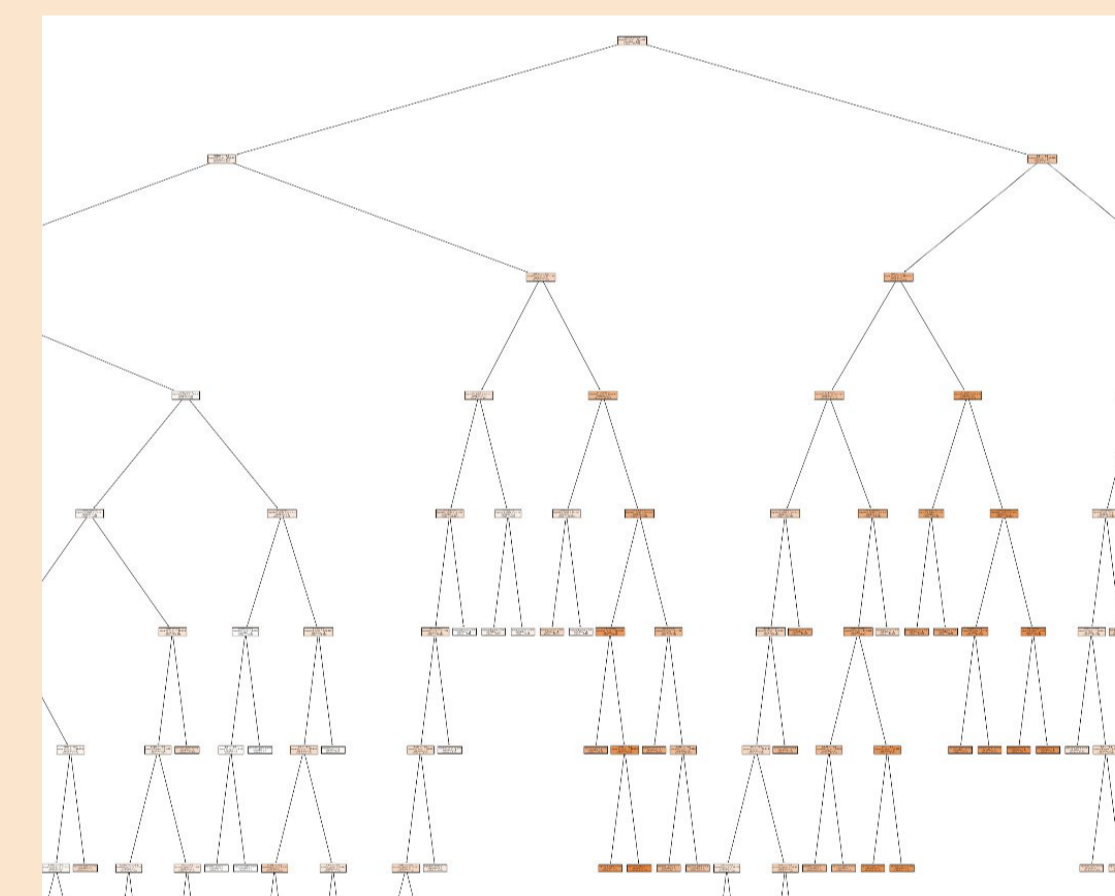
**eXplainable Artificial Intelligence (or XAI)** proposes the creation of ML techniques that:

- Produce more **explainable** models while maintaining **high performance**;
- Enable humans to **understand, trust and manage** the emerging generation of **AI partners**.

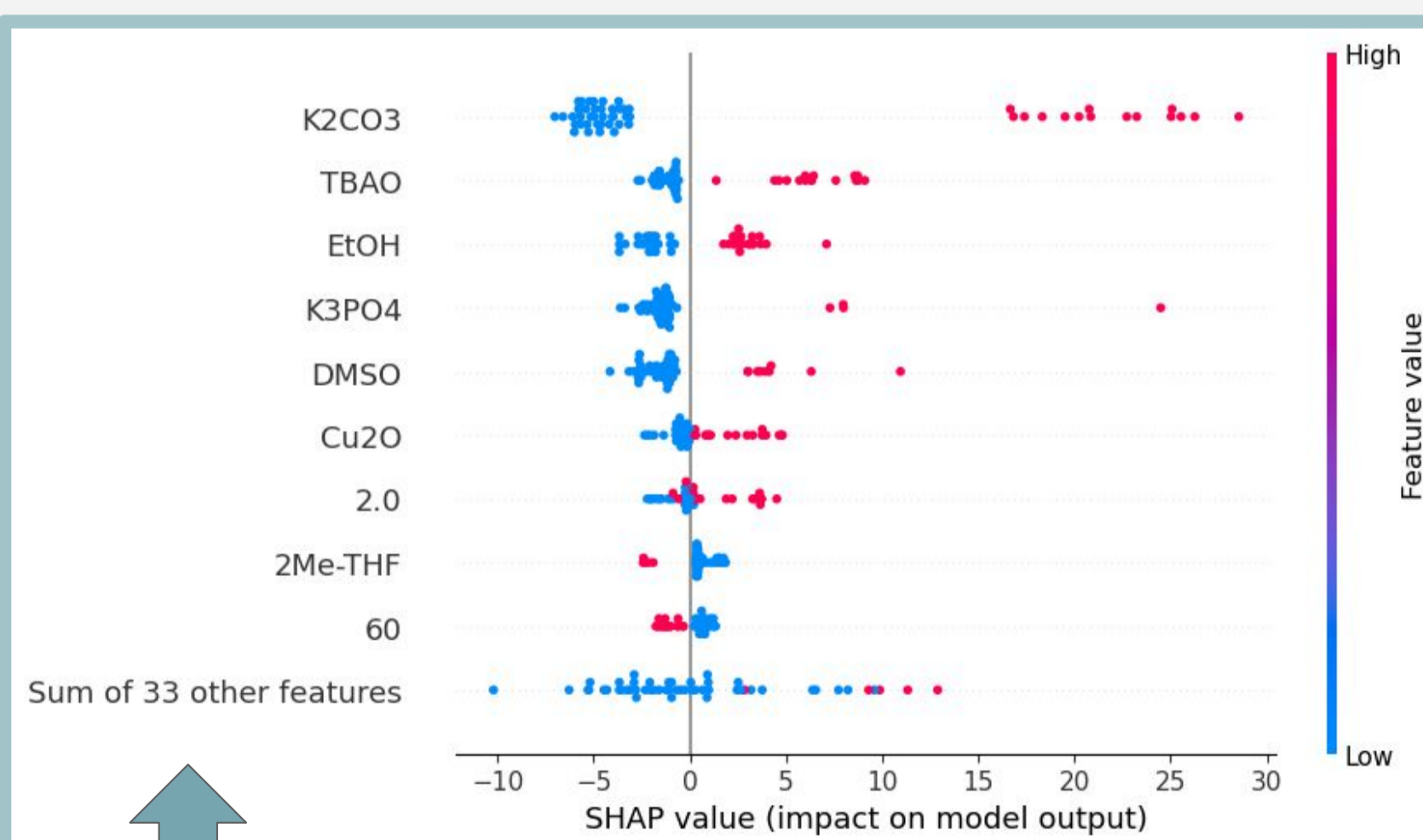
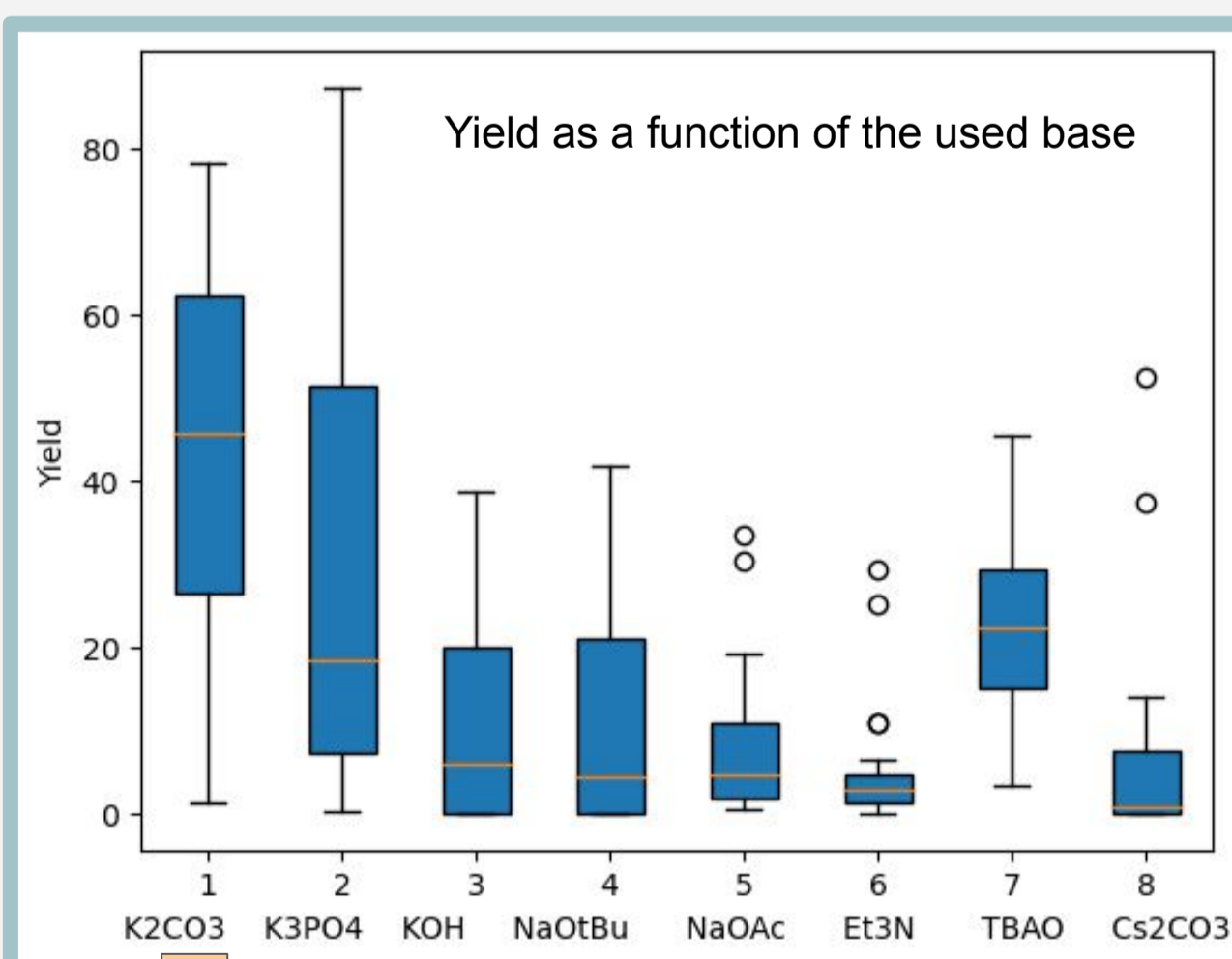
### A select XAI technique : SHapley Additive exPlanations (SHAP)

**1953: Shapley values** by Lloyd Shapley : coalitional game theory method to fairly divide payout among players  
**2017: SHAP values** by Lundberg and Lee prediction of a computer = payout  
feature value of an instance = player

### Decision tree



### From data analysis...



...to feature importance in the model visualized with SHAP.

**This RandomForestRegressor predicted yields with an accuracy of 78% !**

### New reactions ?

By determining which features had the biggest positive impact—biggest positive global SHAP value—on the yield, we can extract the combinations the model deems as most efficient.

The model's propositions are **justified** with SHAP and **align** with the data analysis.

...to natural language explanations of the model's understanding of the data set.

```
shap_global('K2CO3')
```

The mean SHAP value impact on the yield when K2CO3 is present is 21.802705896845396  
The mean SHAP value impact on the yield when K2CO3 is absent is -4.798291795089791  
Overall, the presence of K2CO3 makes the yield higher.

### References:

- Barredo A. et al., Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. Information Fusion, 58:82–115, 2020
- Braconi E., Godineau E. Bayesian optimization as a sustainable strategy for early-stage process development? a case study of cu-catalyzed c–n coupling of sterically hindered pyrazines. ACS Sustainable Chemistry & Engineering, 11(28):10545–10554, 2023
- Scott Lundberg and Su in Lee. A unified approach to interpreting model predictions. arxiv, 2017.

### Conclusion:

By exploiting the data at our disposition, we trained a ML model with **high predictive accuracy** and made explicit its decision process with SHAP. We have shown that **AI is a very promising tool for chemistry** that should keep being experimented with.