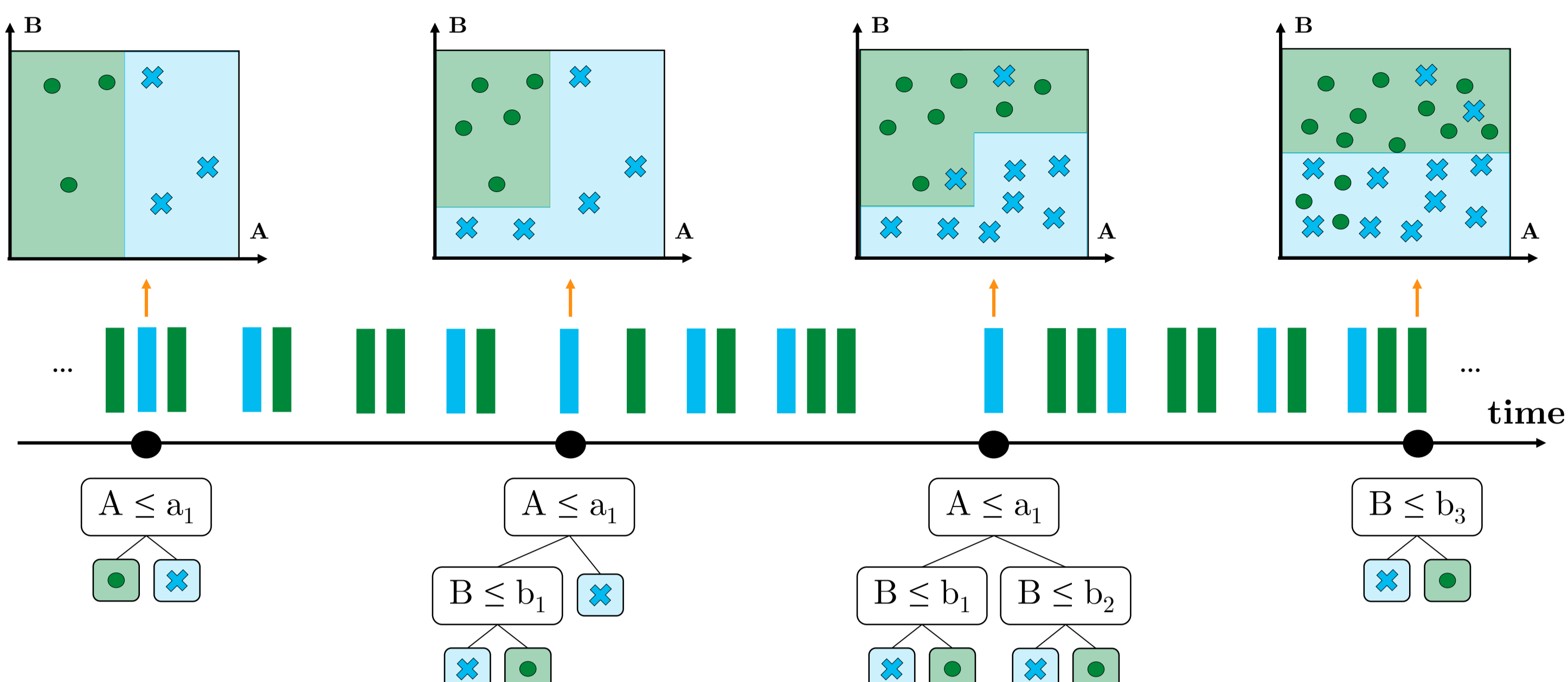


The Problem: Changing Black Box Models

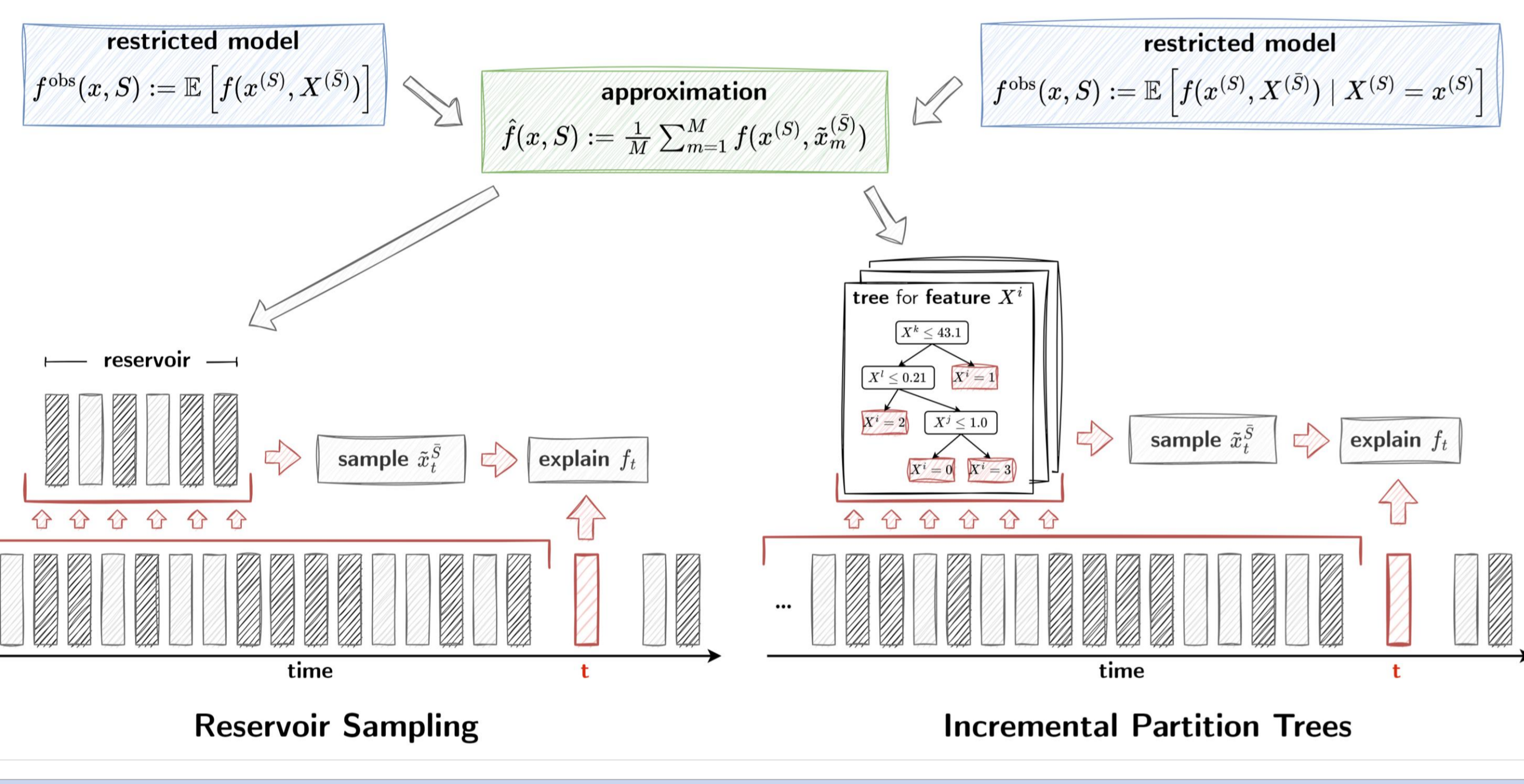


Requires *efficient, any-time* ML models.

Online Feature Removal Mechanisms



Interventional vs Observational



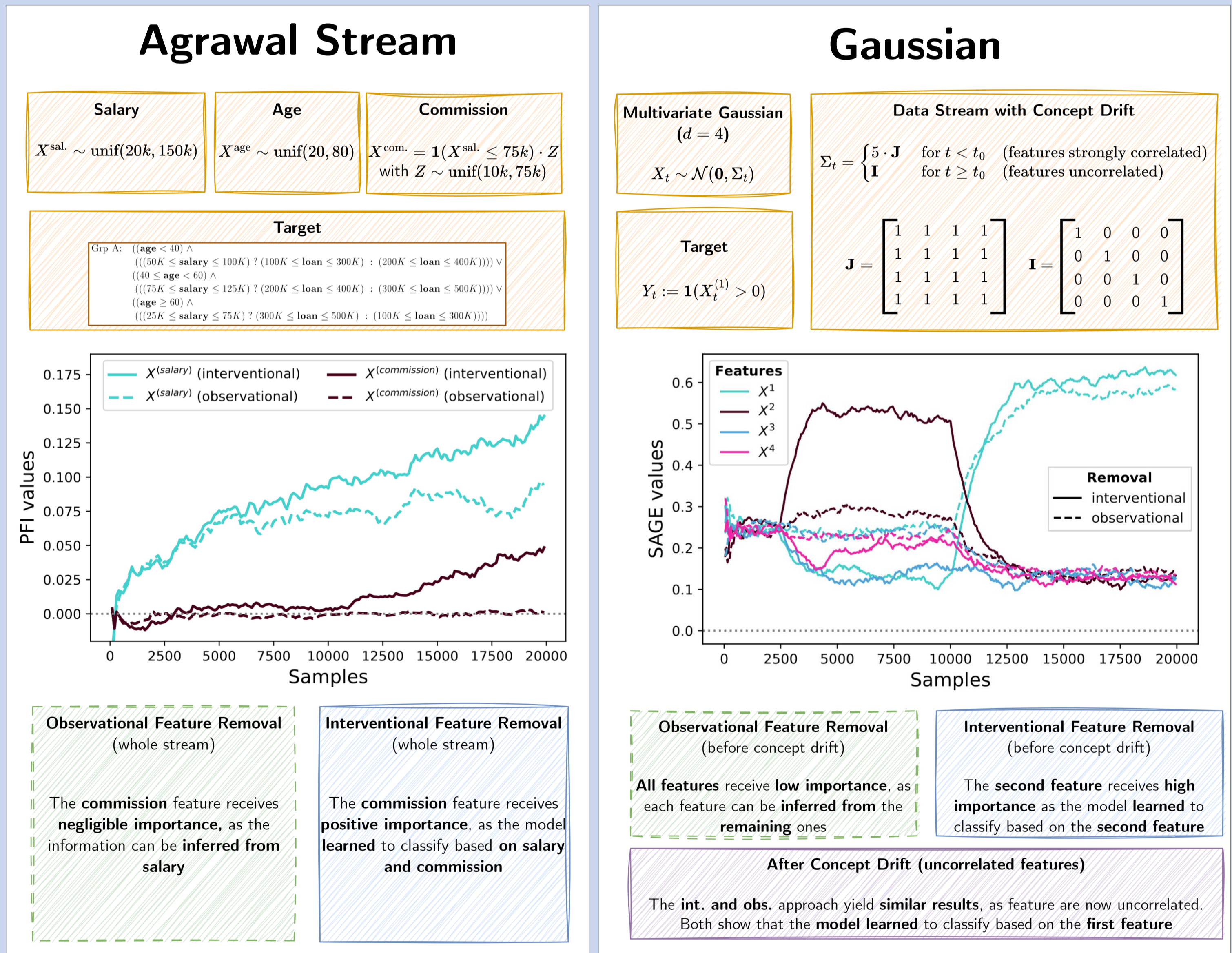
A Solution: Online Explanations

Methods: First incremental, online XAI methods exist.

Global Feature Importance Methods in Dynamic Environments		Global Feature Effect Methods in Dynamic Environments
Incremental Permutation Feature Importance (iPFI) [2]	Incremental Shapley Additive Global Explanation (iSAGE) [3]	Incremental Partial Dependence Plots (iPDPs) [4]
Computes global feature importance incrementally based on the well known permutations tests	Computes global feature importance incrementally based on the Shapley-based SAGE values	Computes global feature effects incrementally based on the established PDPs and ICE curves

Requires *efficient, any-time* feature removal methods.

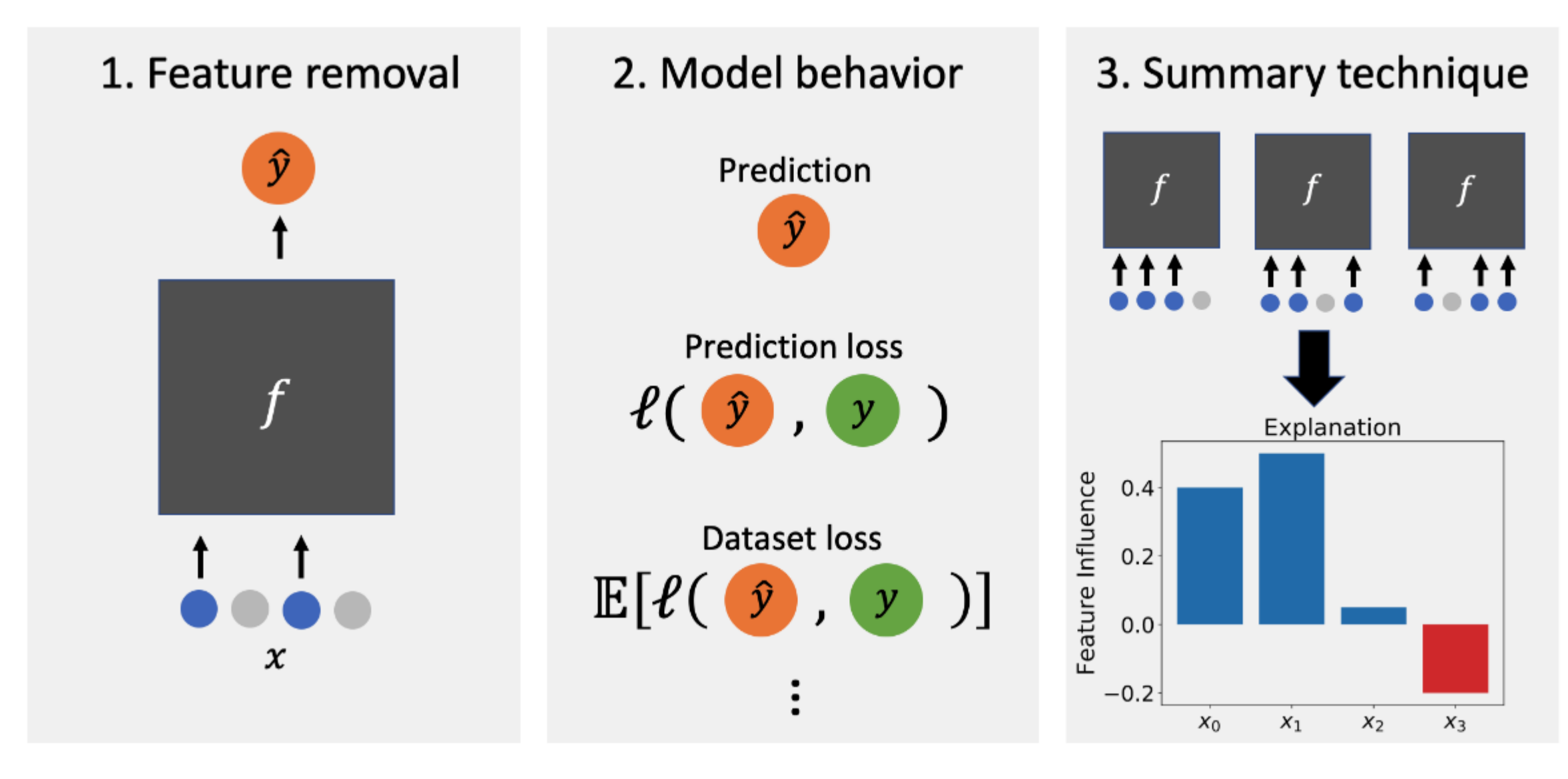
Empirical Comparison



Background: Explaining by Removing [1]

Explaining by removing [1]
 Given a model $f : \mathcal{X} \rightarrow \mathcal{Y}$ on a feature space \mathcal{X} with features $D = \{1, \dots, d\}$, we define

- Feature Removal:** A restricted model $F : \mathcal{X} \times \mathcal{P}(D) \rightarrow \mathcal{Y}$, where features in $\tilde{S} := D \setminus S$ are removed.
- Model Behavior:** A model property $\nu : \mathcal{P}(D) \rightarrow \mathcal{Y}$ that is based on the restricted model.
- Summary Technique:** An aggregation method for different evaluations of ν .



Conclusion

Observational Feature Removal
 Includes the dependencies of features into the explanation. Reveals the information that a feature provides to the model. We confirm that this approach is true to the data [5].

Interventional Feature Removal
 Breaks feature dependencies when computing the explanation. Reveals more accurately what the model has learned. We confirm that this approach is true to the model [5].

Future Work

- Human-grounded evaluation: Conduct user studies and investigate how to explain dynamic models **efficiently** and **effectively**.
- Move to Higher Orders: Move from feature importance to feature interactions.

Open Source: iXAI

docs passing | py310 v0.1.3 | status alpha | License MIT

Installation

```
pip install ixai
```

Quickstart

```
>>> for (n, (x, y)) in enumerate(stream, start=1)
...     accuracy.update(y, model.predict_one(x)) # inference
...     incremental_pfi.explain_one(x, y) # explaining
...     model.learn_one(x, y) # learning
```

- iXAI currently contains four explanation methods: iPFI [2], iSAGE [3], iPDP [4], and MDI [7].
- iXAI allows for **interventional** and **observational** feature removal by means of different imputer systems.
- Looking for **collaborators**!