

# Privacy, Explainable AI

Francesca Naretto, Anna Monreale

University of Pisa

[francesca.naretto@unipi.it](mailto:francesca.naretto@unipi.it)



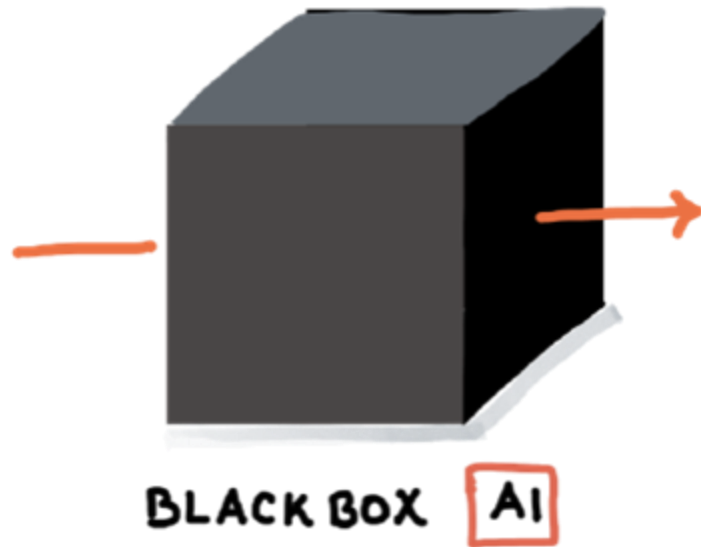
Future  
Artificial  
Intelligence  
Research



**SOBIGDATA**.it  
ITALIAN RESEARCH INFRASTRUCTURE



# Context



We want to explain the global behavior by using ***Global Explainers***

We want to explain the decision on a, instance by using ***Local Explainers***

A ***Machine Learning*** model whose internals are either ***unknown*** to the observer or they are known but ***uninterpretable*** by humans.

# Potential Risk

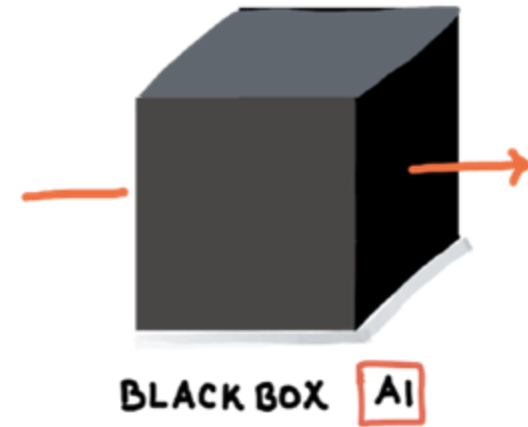
Machine Learning models may enable privacy risks

There are several privacy attacks design to attack Machine Learning models, such as:

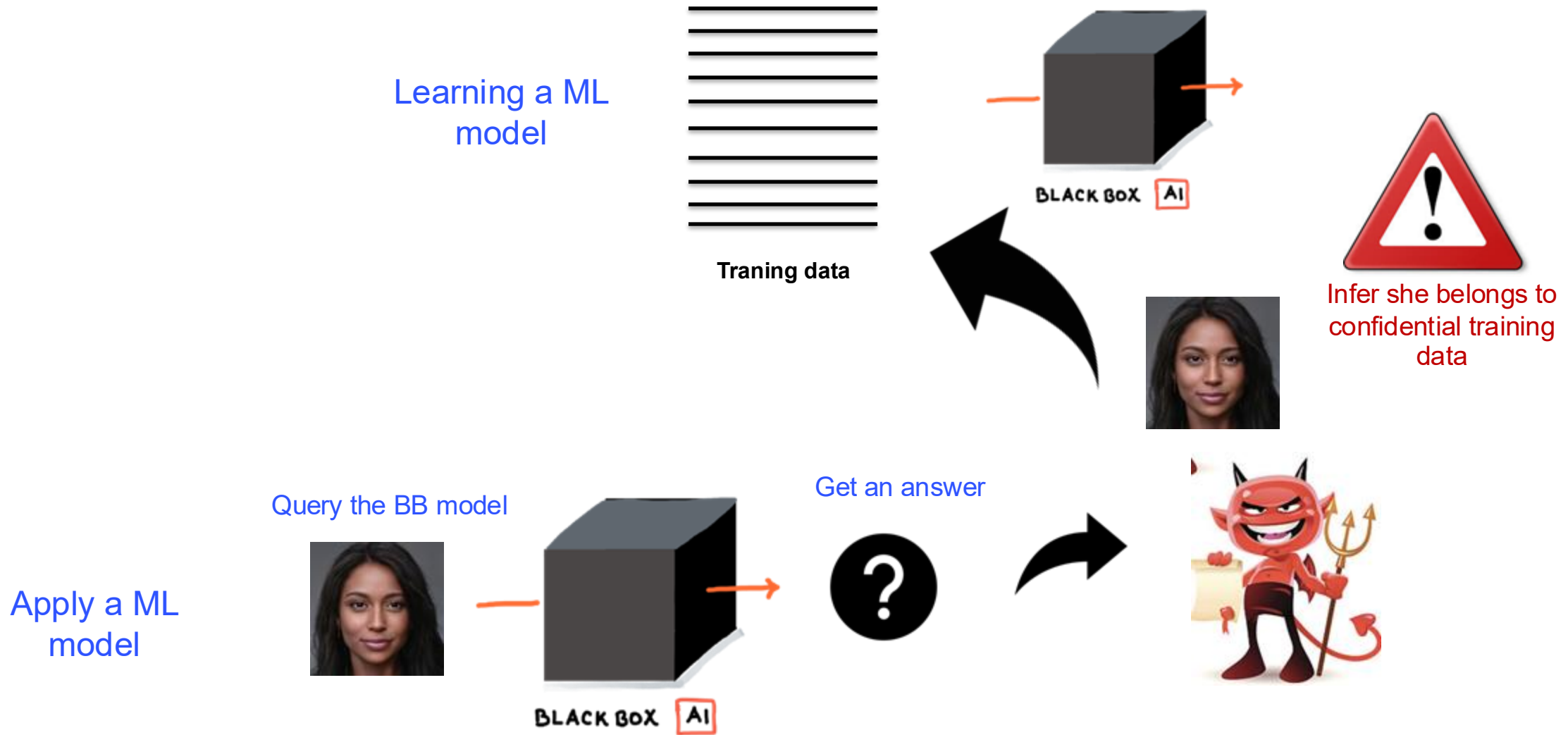
1. Membership Inference attack
2. Reconstruction attack
3. Property inference attack

Explainers are themselves Machine Learning models, even if interpretable!

So, they can also be attacked.

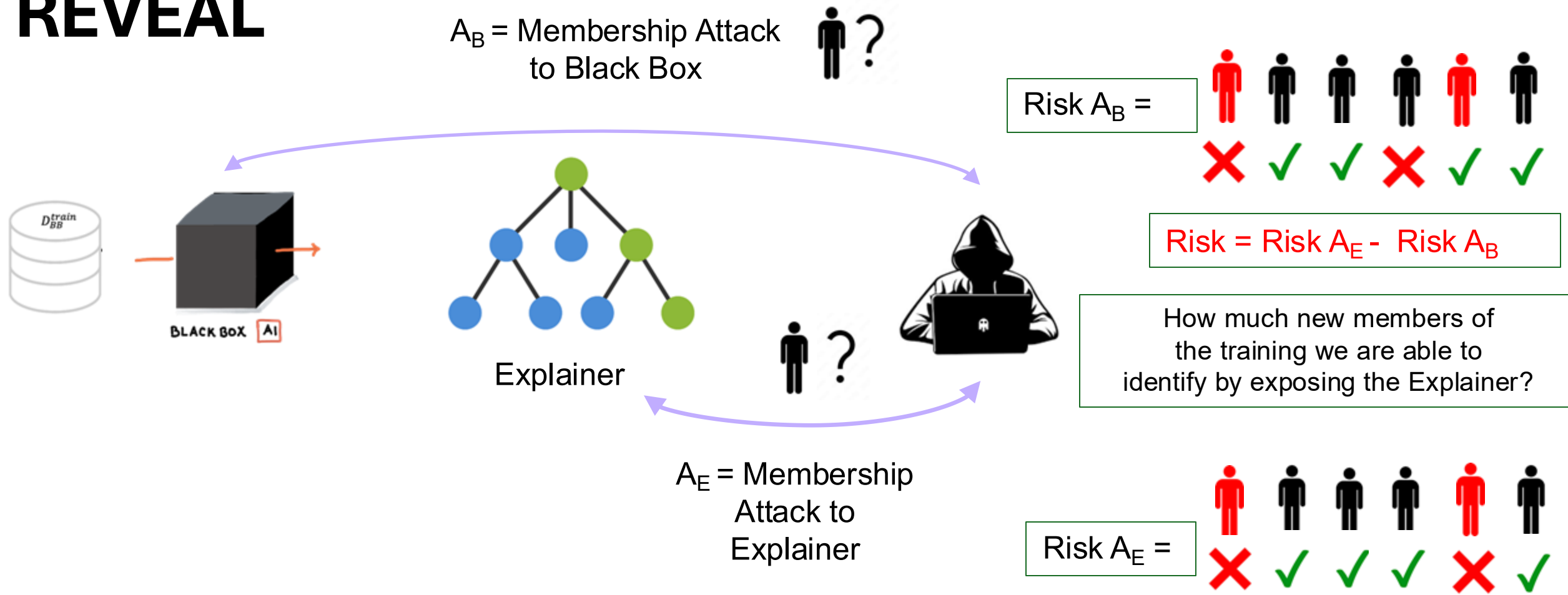


# Privacy Risk Assessment

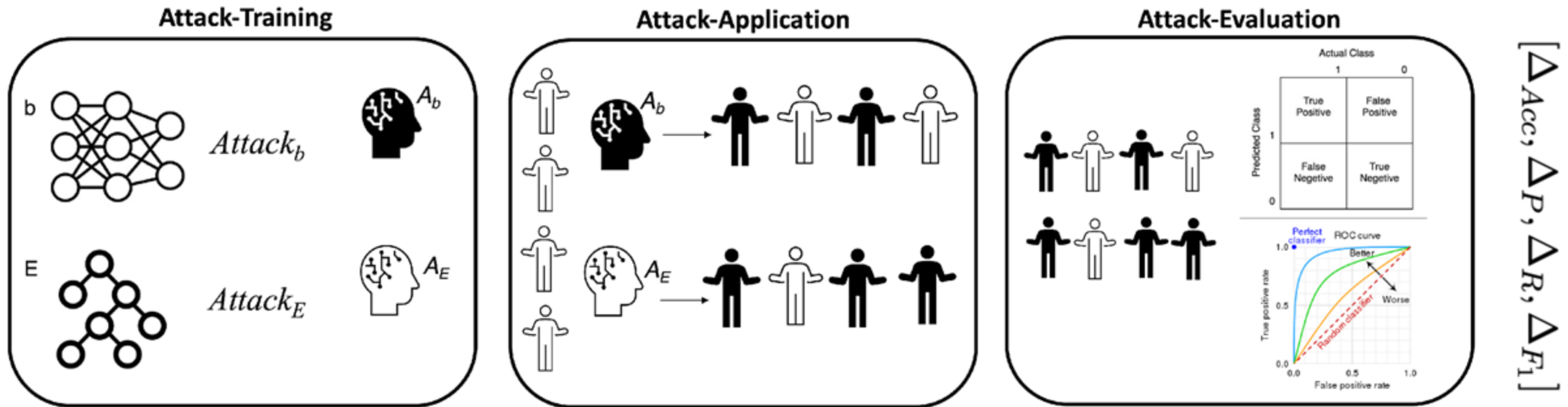


# What about Explainers and privacy risks?

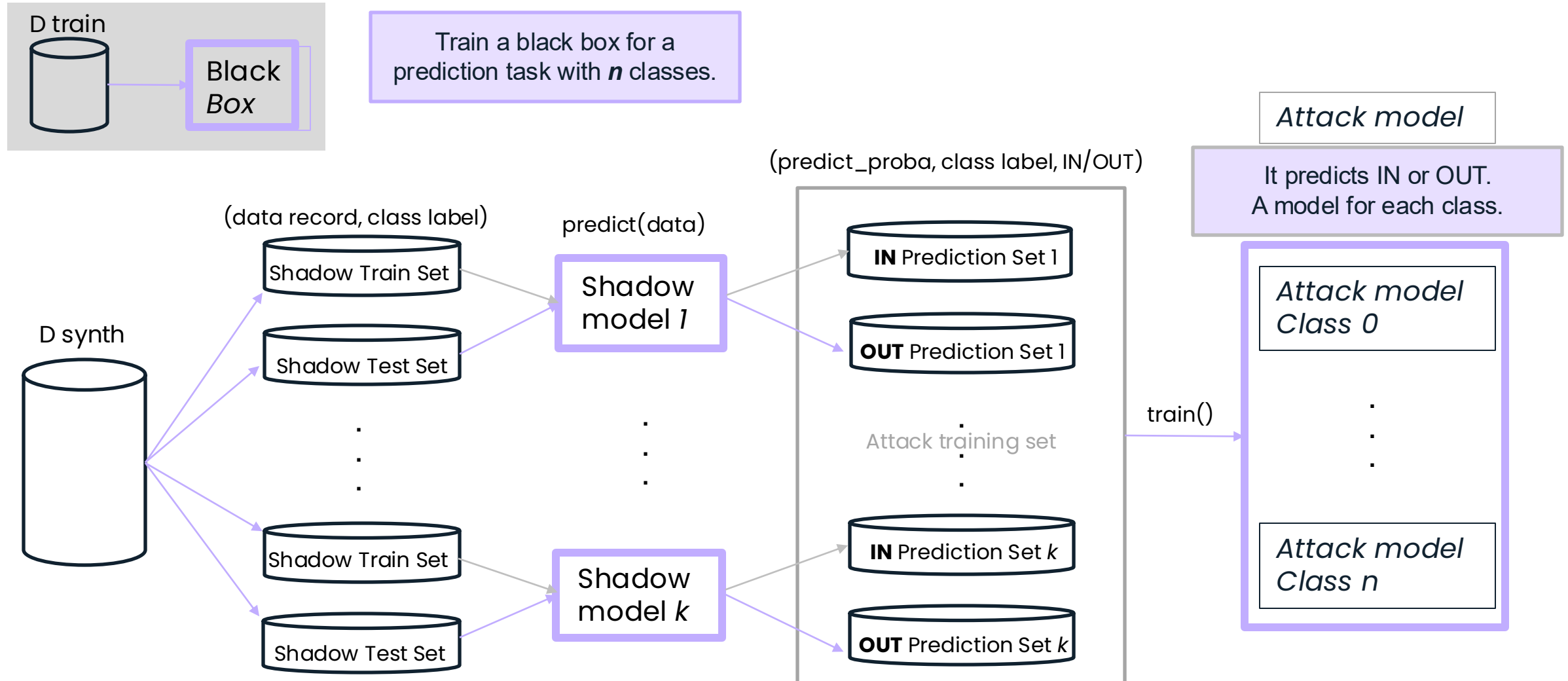
## REVEAL



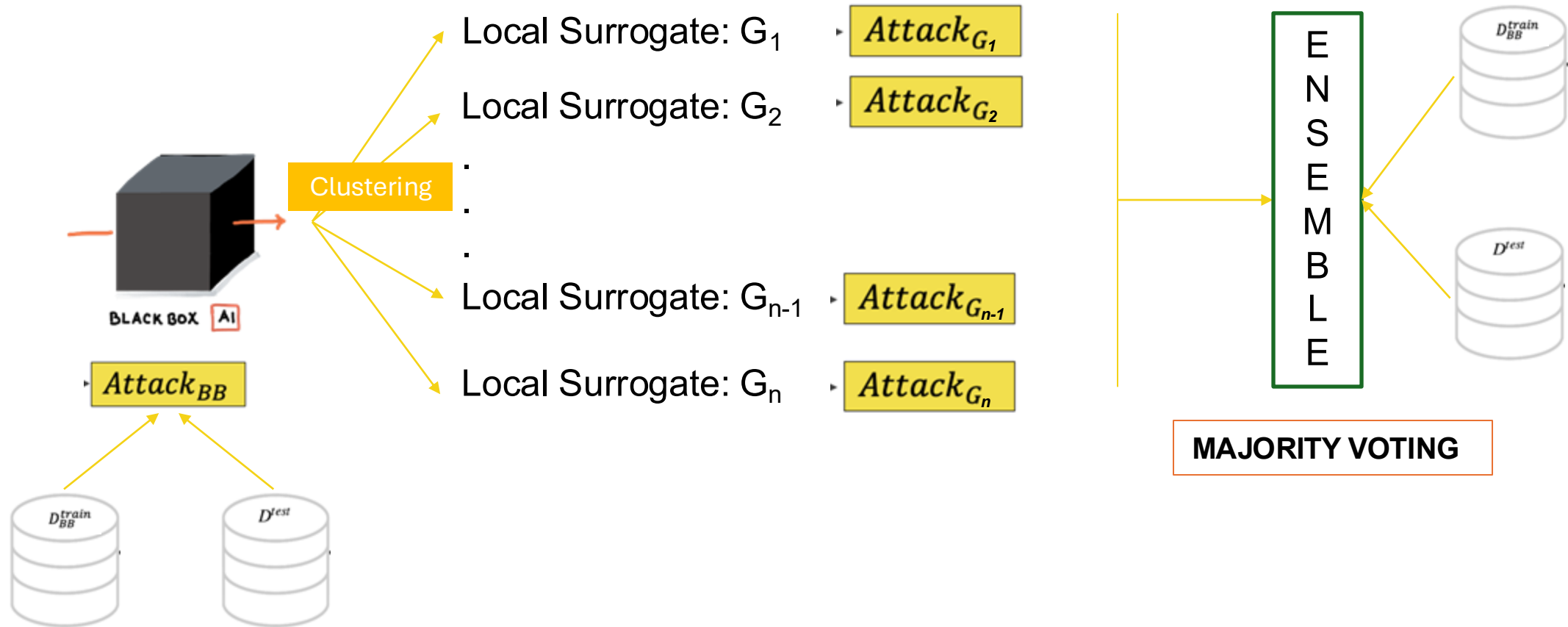
# REVEAL: privacy exposure of explainers



# The privacy attack: MIA



# What about Local Explainers?





# **FastSHAP<sup>++</sup>**

## **A federated private explainer end-to-end**

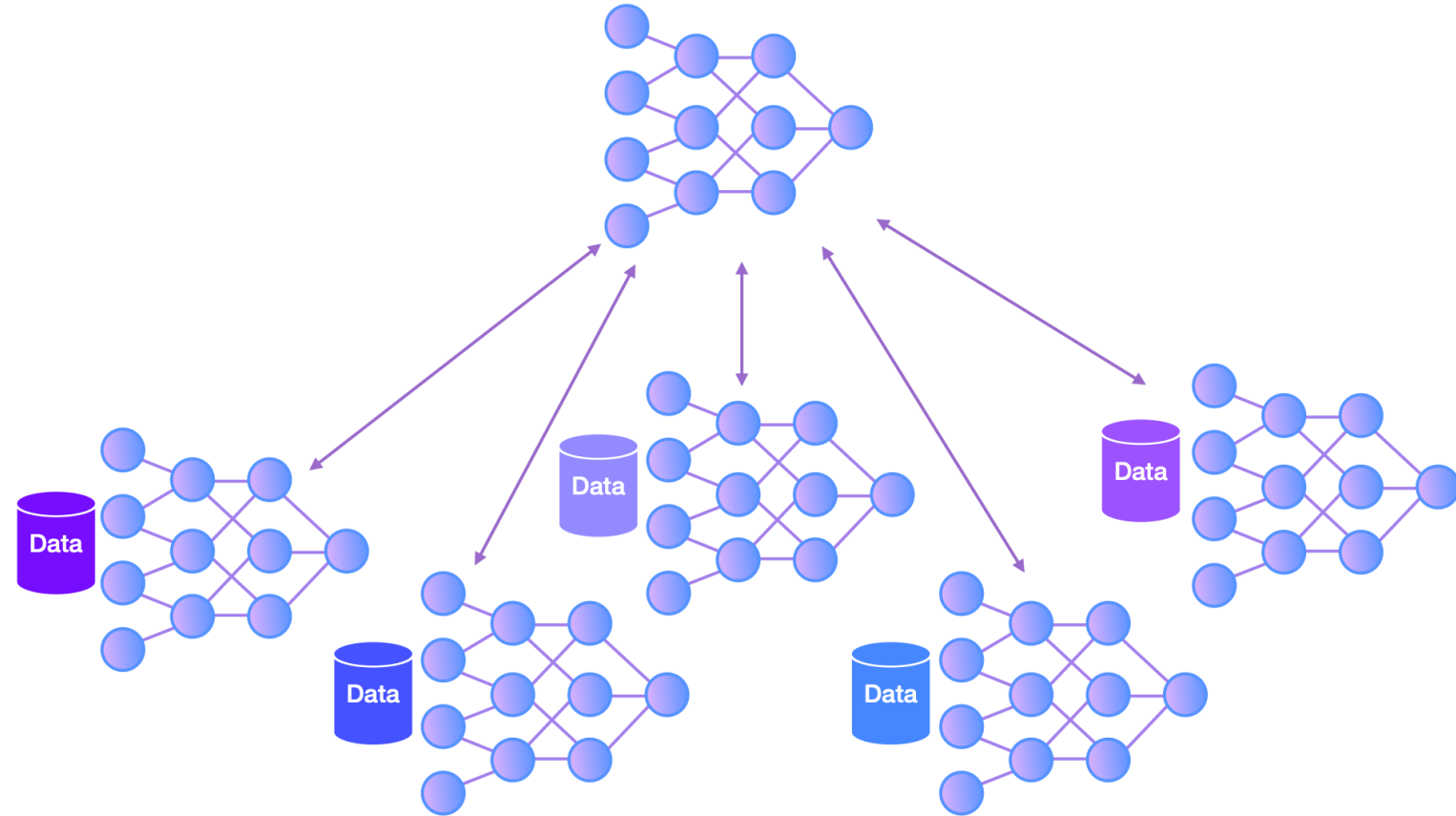
Valerio Bonsignori, Luca Corbucci, *Francesca Naretto*, Anna Monreale

Is it possible to explain Federated Learning models while preserving privacy and Federated Learning constraints?

Is it possible to explain **Federated Learning** models while preserving privacy and Federated Learning constraints?

# Federated Learning

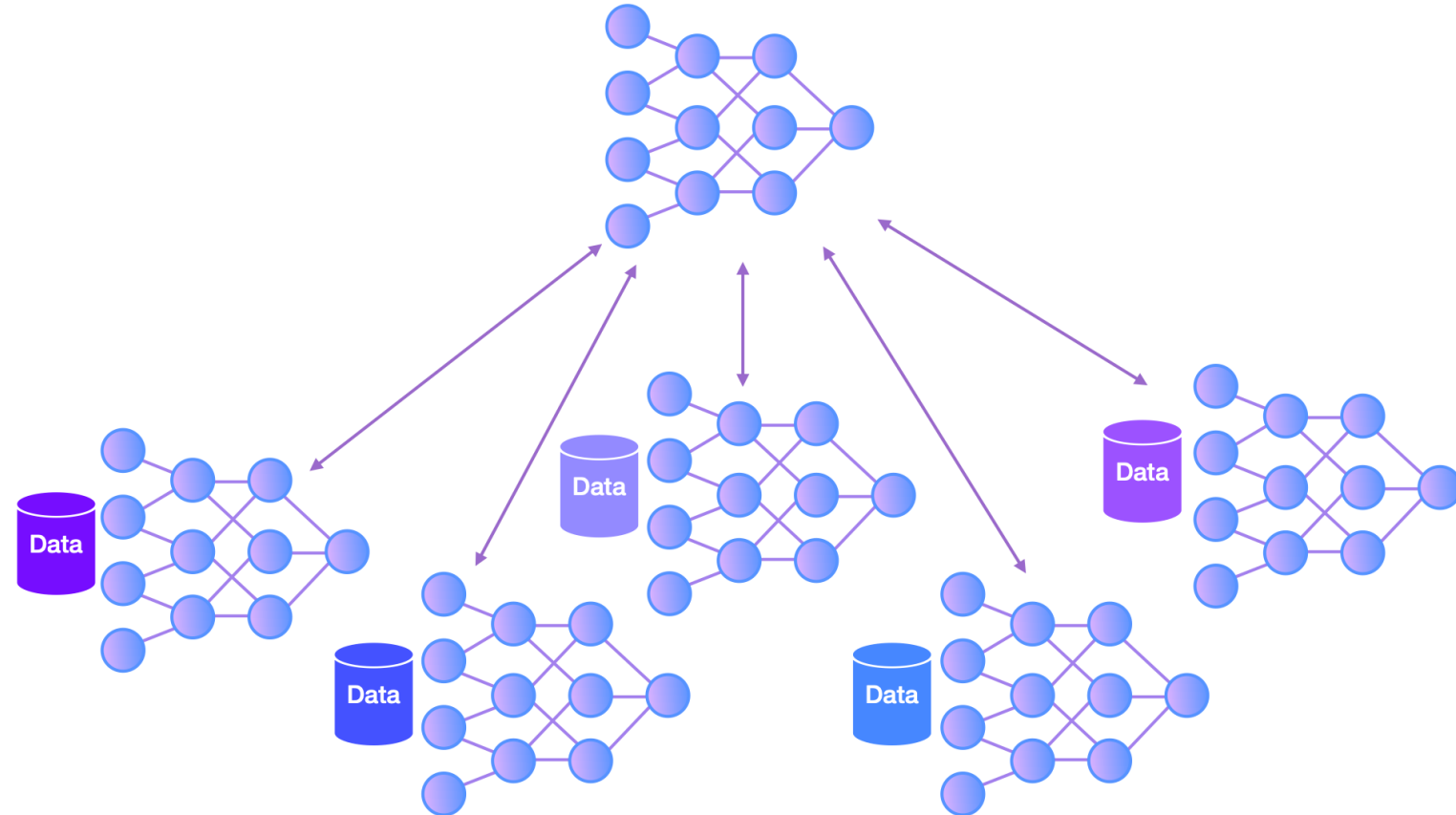
Clients don't share their data, they only exchange model updates.



# Federated Learning

## FedAvg

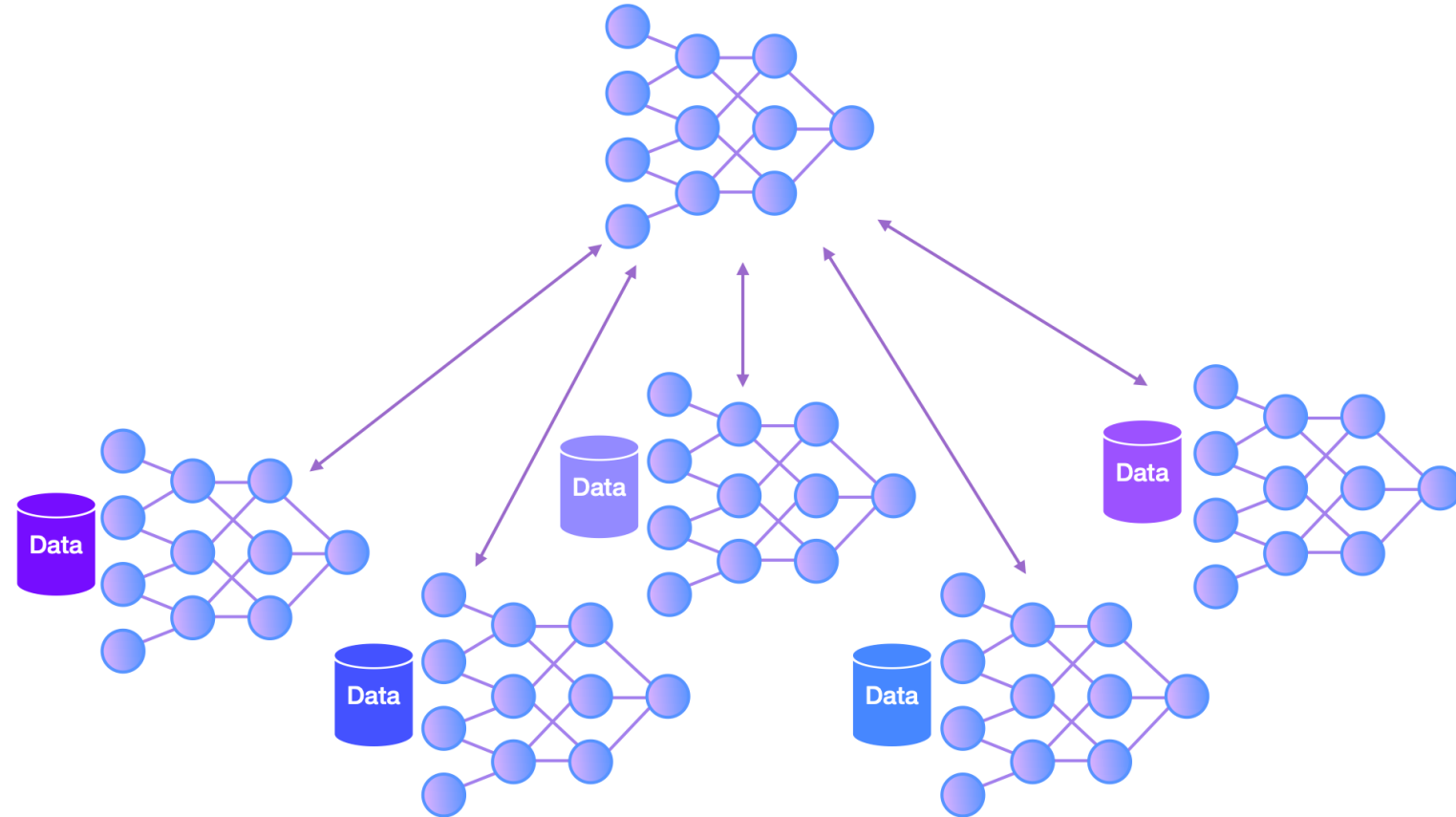
The server updates the global model by computing an average of the local parameter vectors returned by the participating clients after their local optimization steps.



# Federated Learning

Clients don't share their data, they only exchange model updates.

Good generalization capabilities.



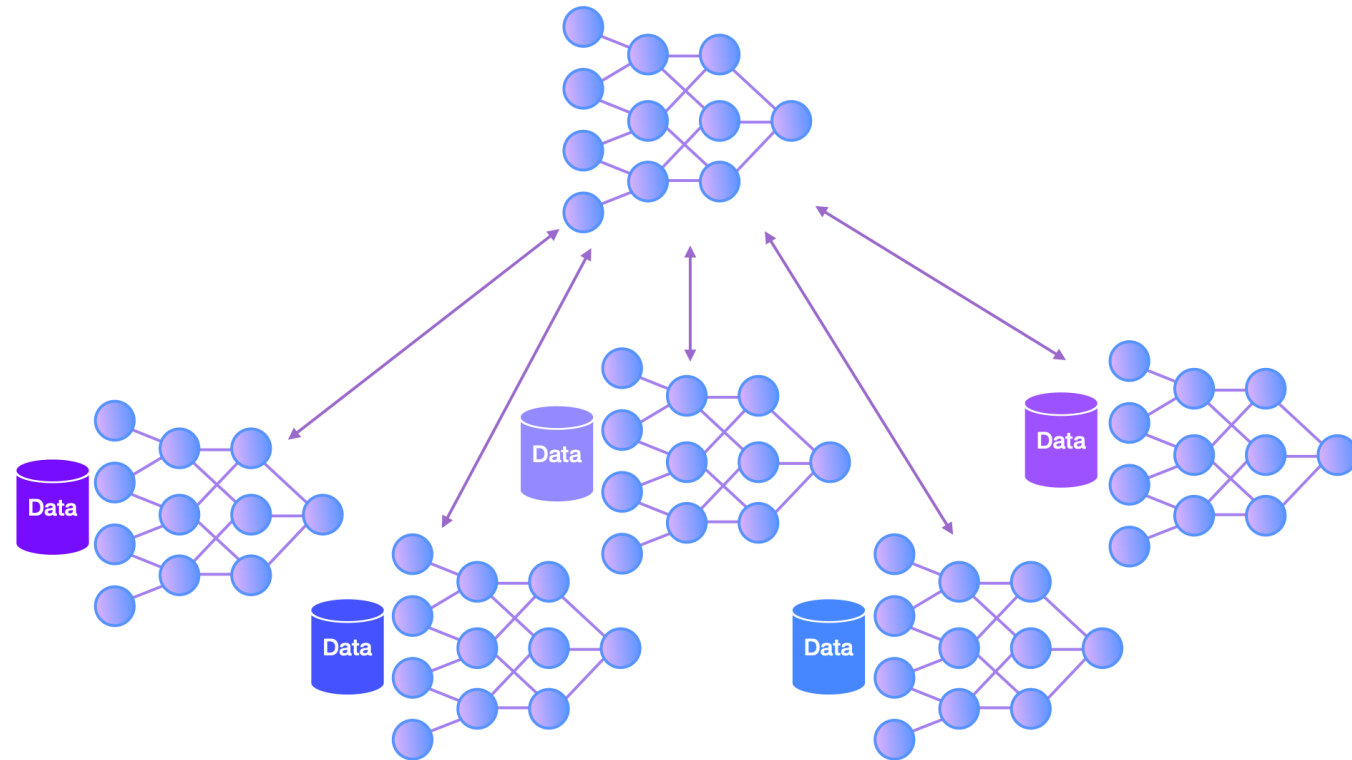
Is it possible to explain Federated Learning models while preserving **privacy** and Federated Learning constraints?

# Federated Learning & Privacy

Privacy attacks are still possible.

- Inverting gradients attacks
- Membership Inference Attacks
- Property Inference Attacks

Differential private learning of ML model can support privacy protection

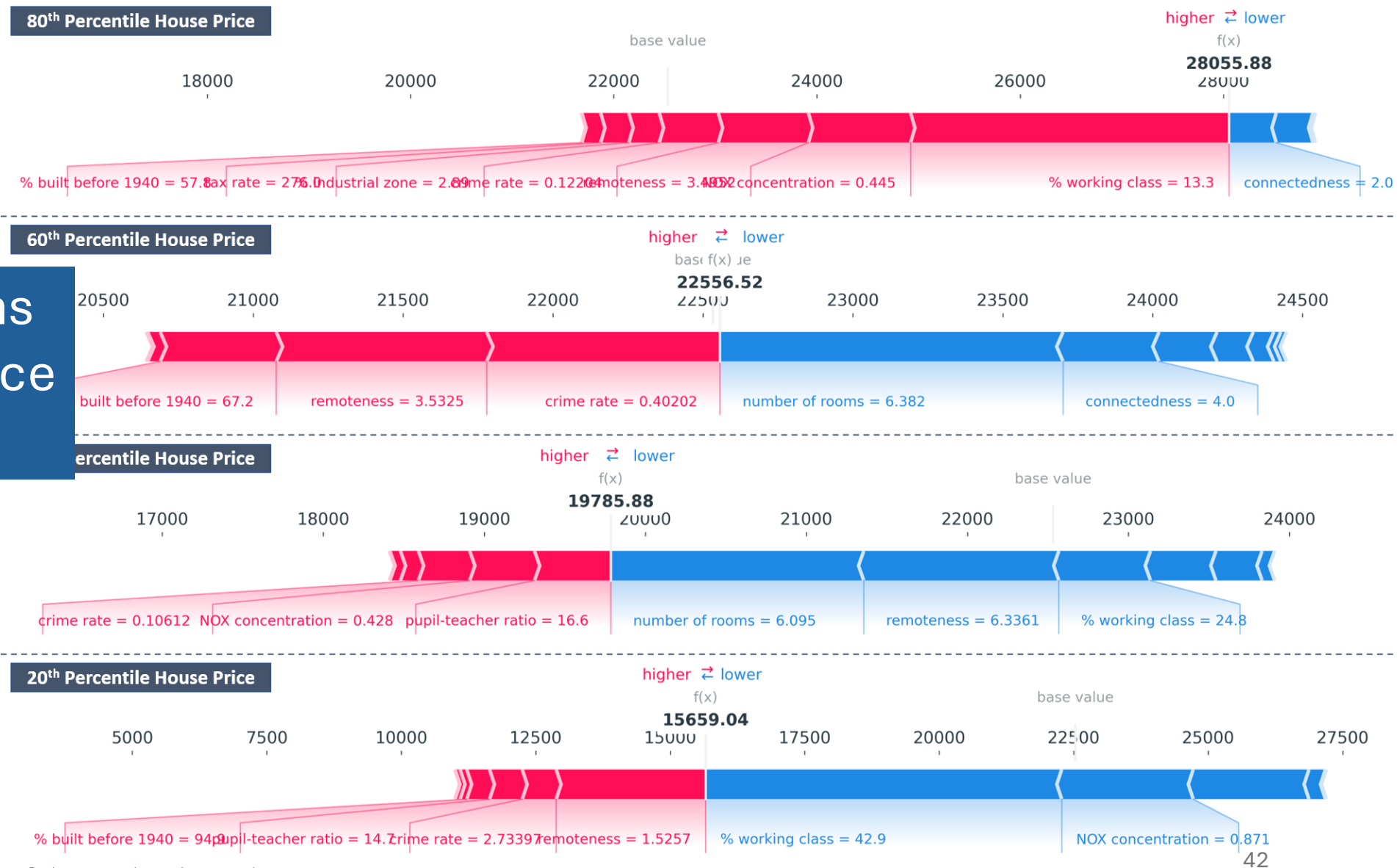




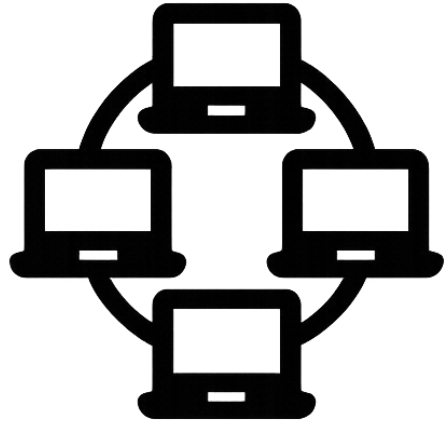
Is it possible to **explain** Federated Learning models while preserving privacy and Federated Learning constraints?

# Explainable AI - SHAP

- Local explanations
- Feature importance
- Agnostic



# We would like to:



Work in a Federated Learning scenario

Have local explanations

Preserve the privacy during all the steps of the pipeline

# Limits



SHAP can be slow



Federated Learning VS explanations

- SHAP requires data to be trained on
- We don't have them on the server side
- Problems with privacy



Overall... Limited privacy protection

A possible solution is to use **FastSHAP**

- An explainer Neural Network
- Principles of SHAP values are still respected
- Good trade off between accuracy and speed up



# FastSHAP

A possible solution is to use **FastSHAP**

- Not tailored for Federated Learning settings
- No privacy protection
- It uses original data for training



# FastSHAP<sup>++</sup>

## Federated Learning

No exchange of data.

## Fast Explanations

**FastSHAP** Explainer generates explanations in a forward step.

## Local Explanation

Feature Importance Explanations using **FastSHAP**.

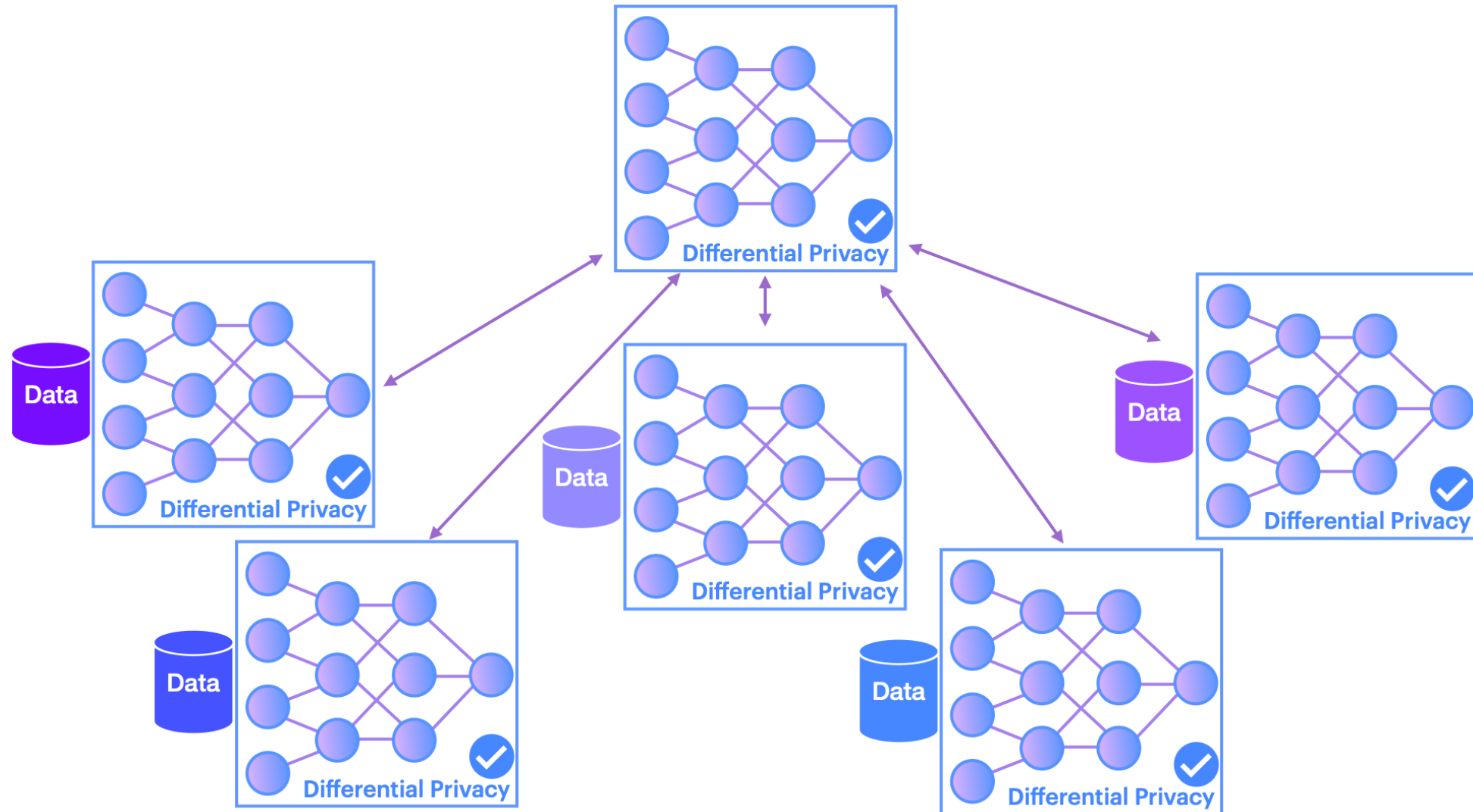
## Distributed Explainer Training

Explainer complies with Federated Learning constraints.

## Privacy Protection

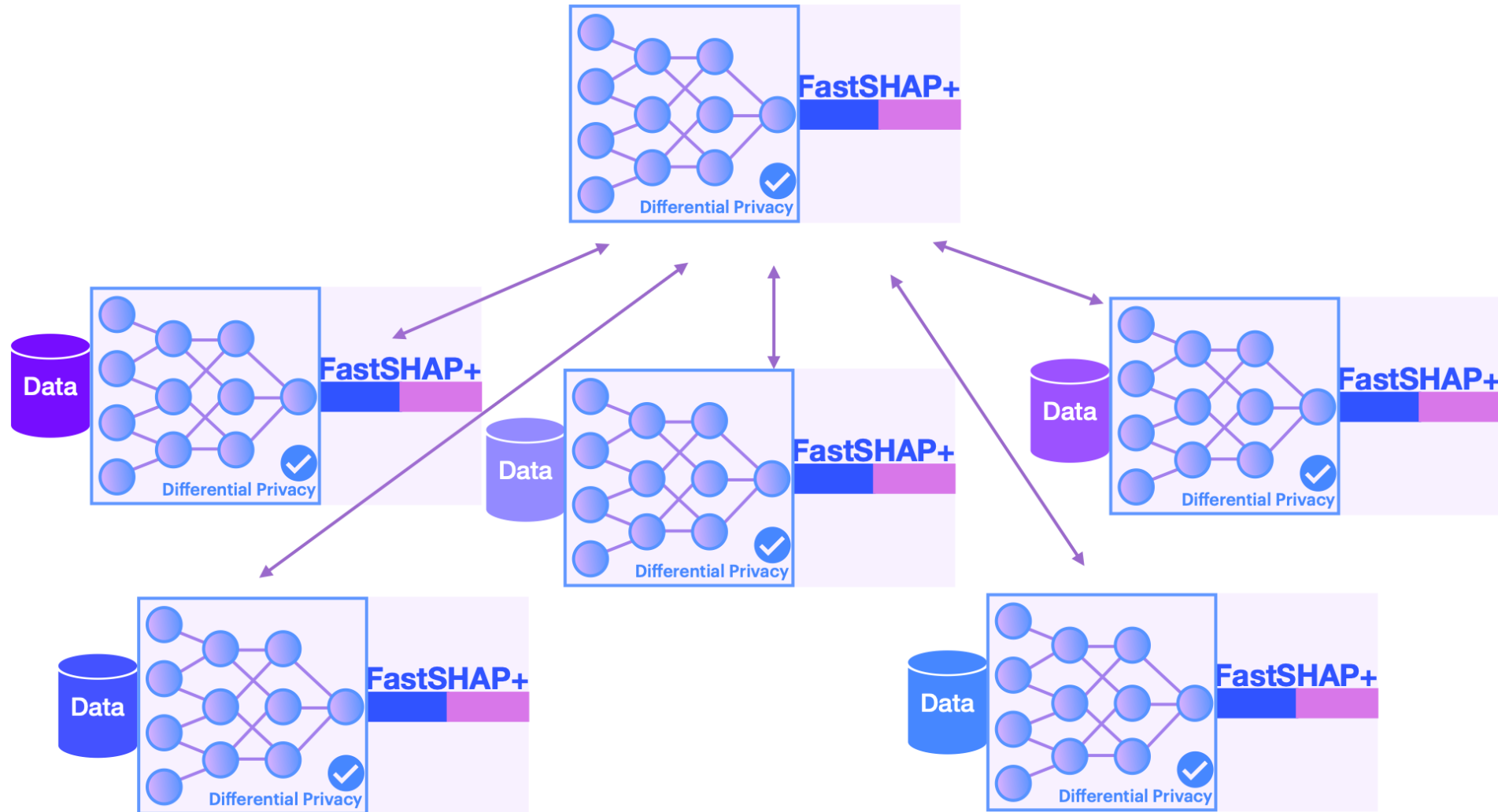
A fully protected pipeline for privacy-guarantee explanations.

# FastSHAP++

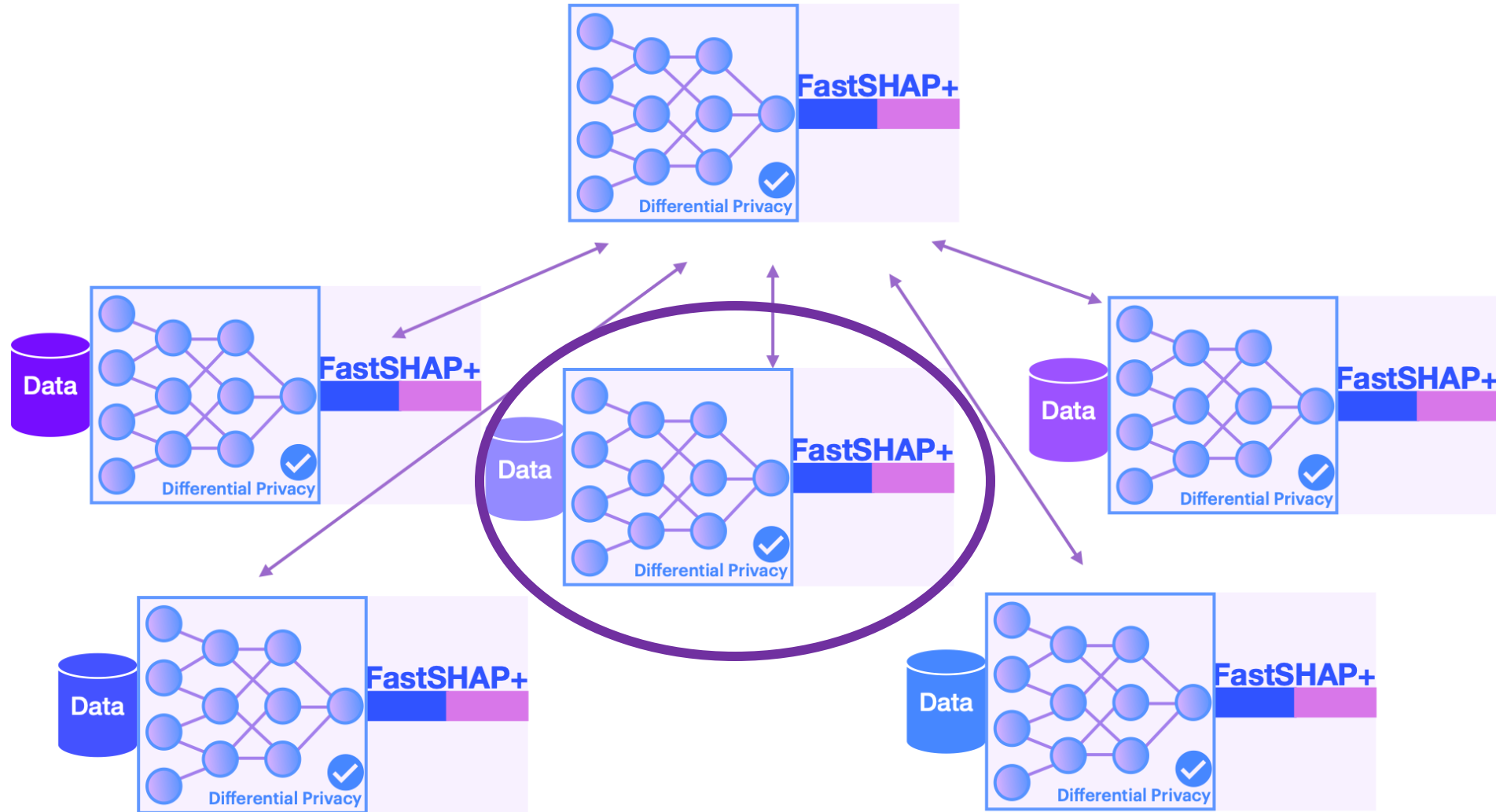




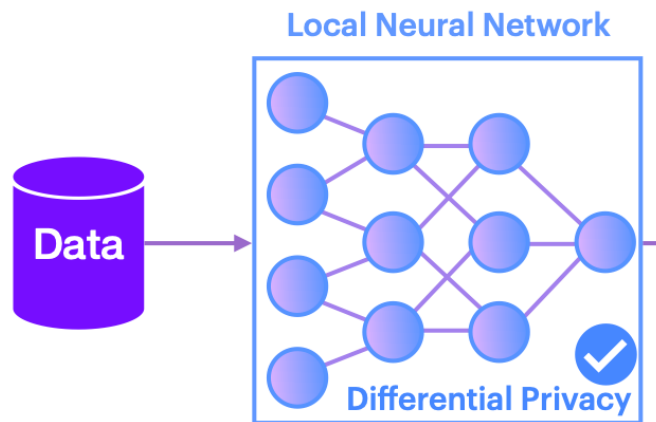
# FastSHAP++



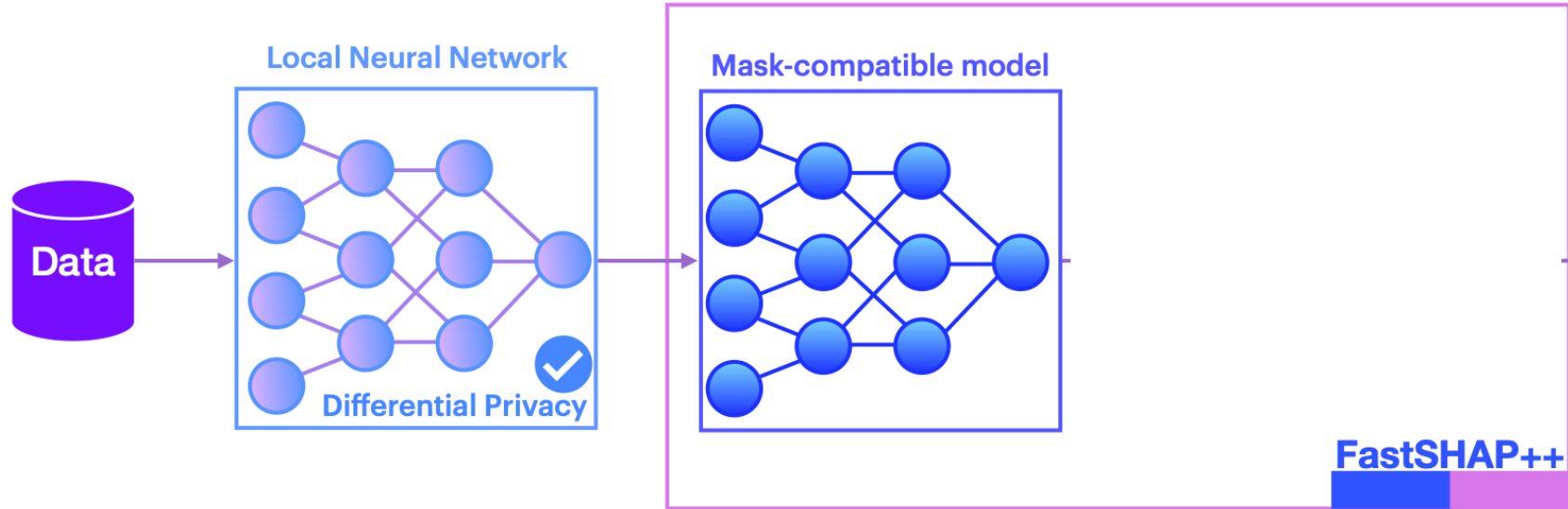
# FastSHAP++



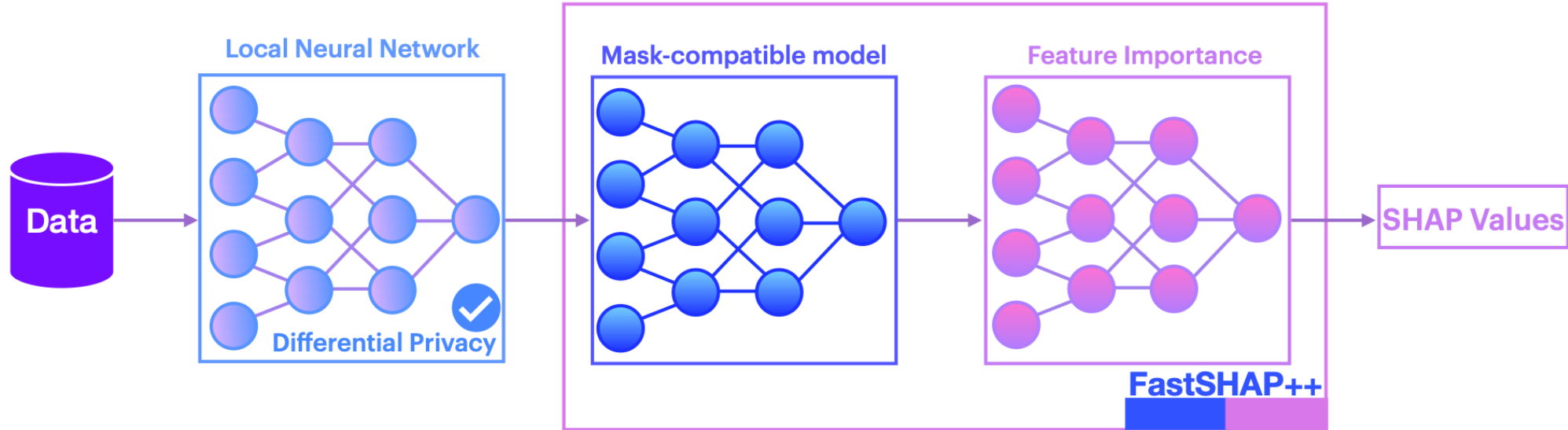
# FastSHAP++



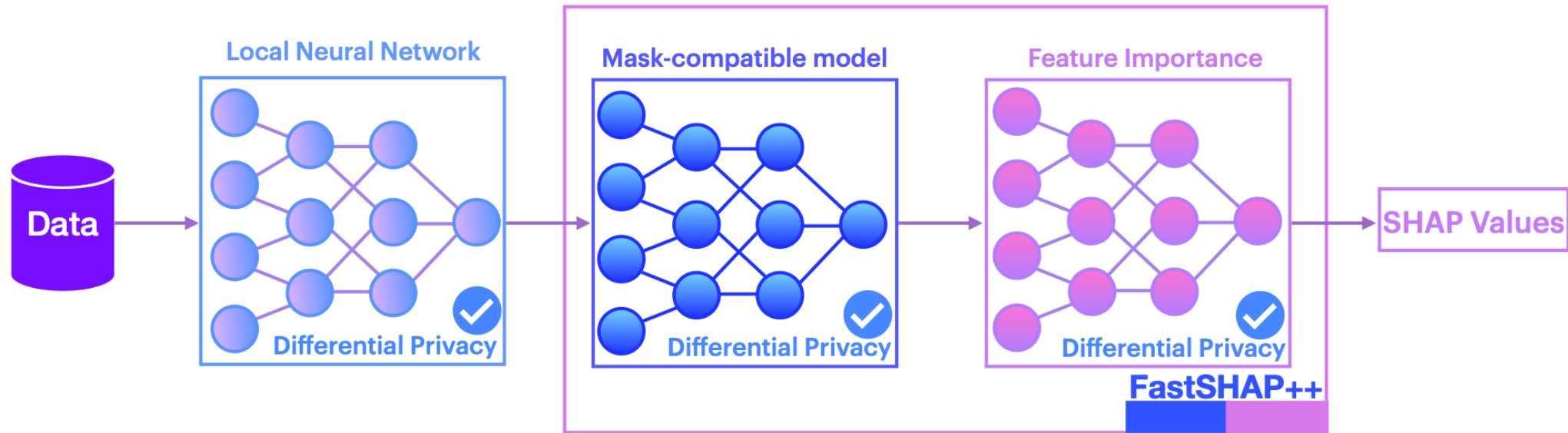
# FastSHAP++



# FastSHAP++

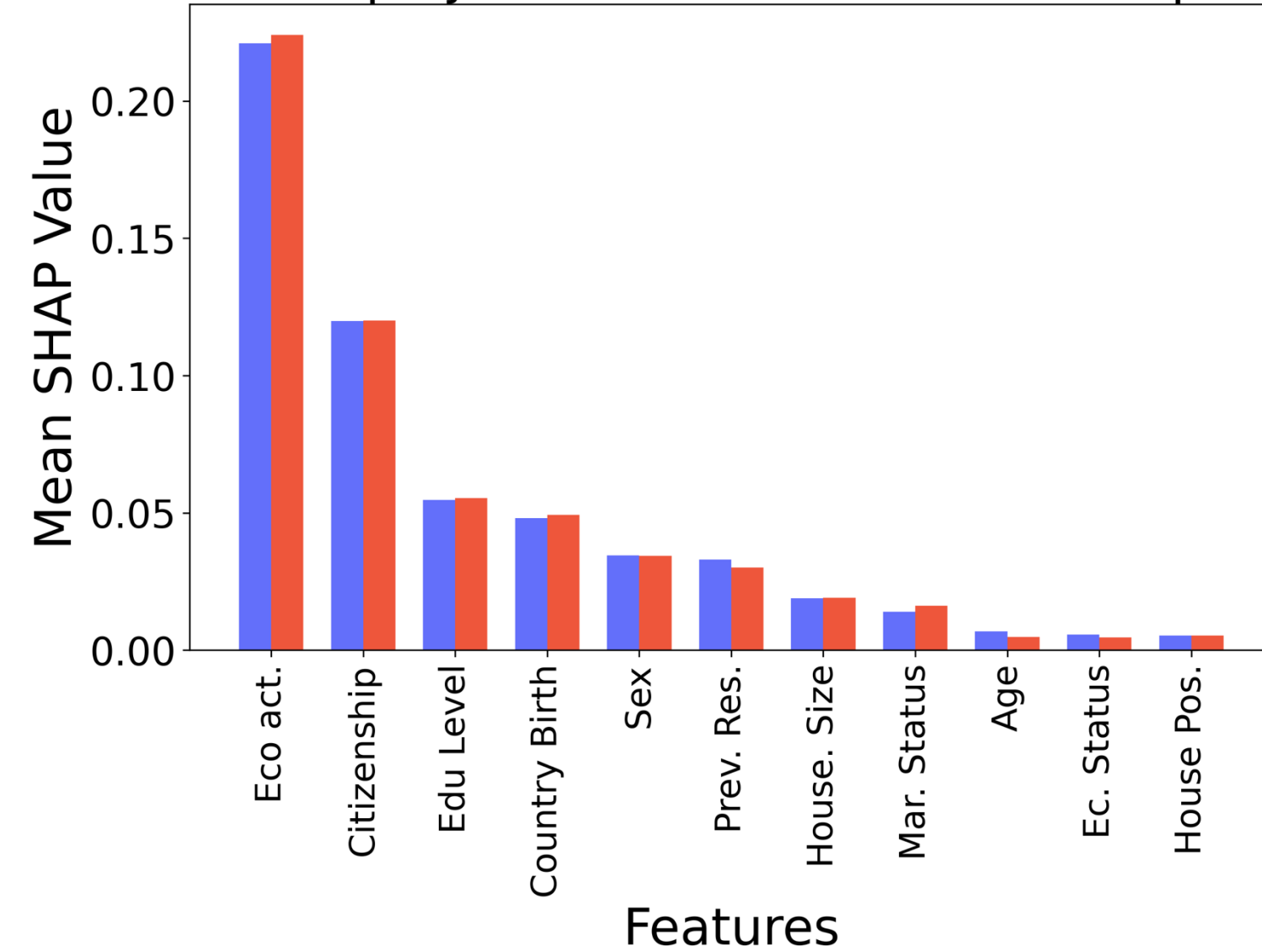


# FastSHAP++



# Are the explanations matching the centralized?

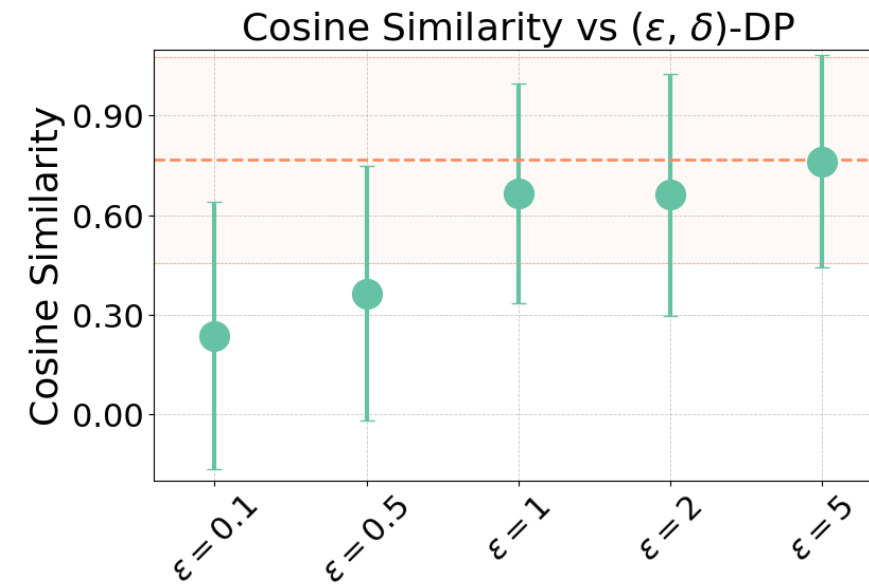
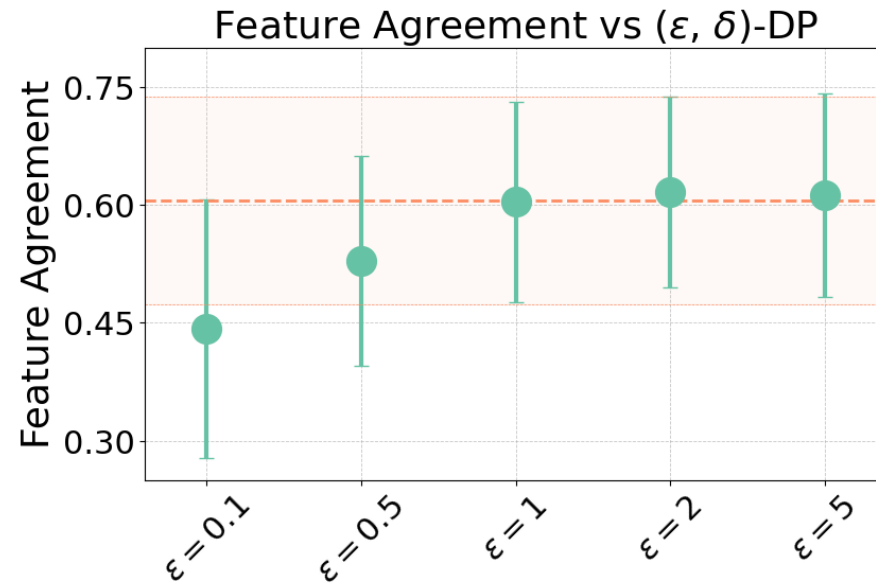
Mean Shapley Value: FL vs Centralized Explainer



## Centralized VS Federated

- Same direction of importance
- Similar magnitudes

# Evaluation on Privacy



FastSHAP++ no privacy VS FastSHAP++ with privacy

The quality of the explanations is good with  $\epsilon \geq 1$ .



# FastSHAP<sup>++</sup>

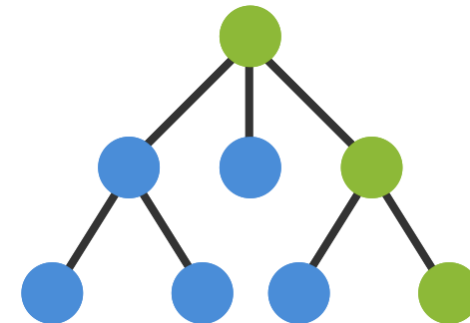
FastSHAP<sup>++</sup> achieves centralized-level explanation quality while preserving clients' data privacy.

It achieves privacy without degrading the explanation quality too much, particularly when  $\varepsilon \geq 1$ .

# Global Explainer: TREPAN

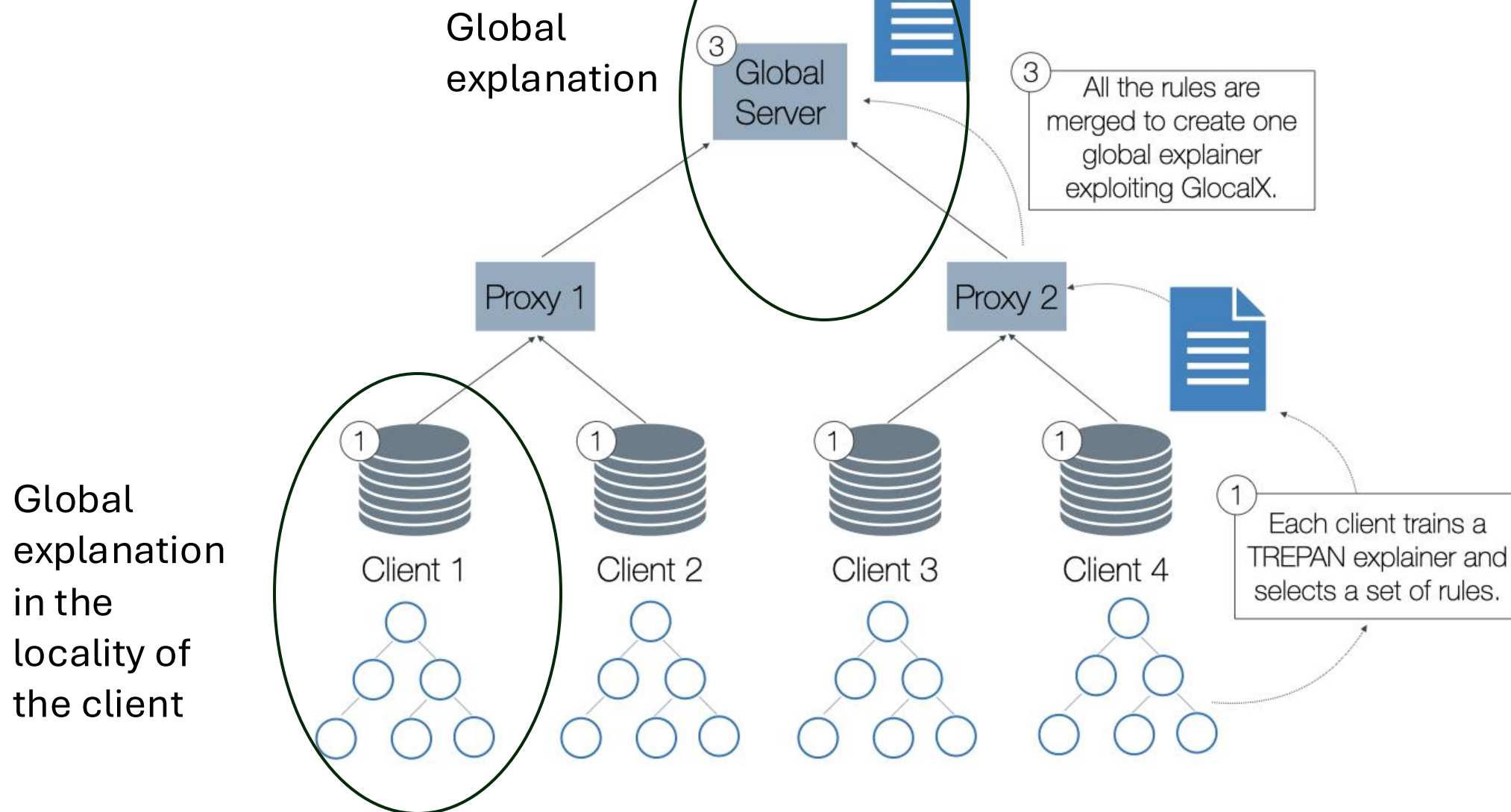
```
1 T = root_of_the_tree()
2 Q = <T, X, {}>
3 while Q not empty & size(T) < limit
4     N, XN, CN = pop(Q)
5     ZN = random(XN, CN)
6     yZ = b(Z), y = b(XN)
7     if same_class(y ∪ yZ)
8         continue
9     S = best_split(XN ∪ ZN, y ∪ yZ)
10    S1 = best_m-of-n_split(S)
11    N = update_with_split(N, S1)
12    for each condition c in S1
13        C = new_child_of(N)
14        CC = CN ∪ {c}
15        XC = select_with_constraints(XN, CN)
16        put(Q, <C, XC, CC>)
```

- Enriches the training data
- Labels the data by using the BB
- Train a DT (surrogate model)



# GLOB-FLEX

Define a rule-based explainer  
for FL models privacy-  
protected.



# How to merge rules

---

**Input:**  $\mathbb{E}$  explanation theories,  $\alpha$  filter threshold

**Output:**  $E$  explanation theory

```
1:  $E \leftarrow \emptyset$ 
2: repeat
3:    $\mathbb{Q} \leftarrow \text{SORT}(\mathbb{E})$                                  $\triangleright$  sort pairs of theories by similarity
4:    $\text{merged} \leftarrow \text{False}$ 
5:    $X' \leftarrow \text{batch}(X)$ 
6:   while  $\neg \text{merged} \wedge \mathbb{Q} \neq \emptyset$  do
7:      $E_i, E_j \leftarrow \text{POP}(\mathbb{Q})$                                  $\triangleright$  select most similar theories
8:      $E_{i+j} \leftarrow \text{MERGE}(E_i, E_j, X')$                      $\triangleright$  merge theories
9:     if  $\text{BIC}(E_{i+j}) \leq \text{BIC}(E_i \cup E_j)$  then                 $\triangleright$  verify improvement
10:       $\text{merged} \leftarrow \text{True}$ 
11:      break
12:   if merged then                                             $\triangleright$  merge occurred
13:      $\mathbb{E} \leftarrow \text{UPDATE}(E_i, E_j, E_{i+j})$                  $\triangleright$  update hierarchy
14: until  $|\mathbb{E}| > 1 \wedge \text{merged}$                                  $\triangleright$  until the merge is successful
15:  $E \leftarrow \text{FILTER}(E, \alpha)$                                  $\triangleright$  Filter final theory
16: return  $E$ 
```

---

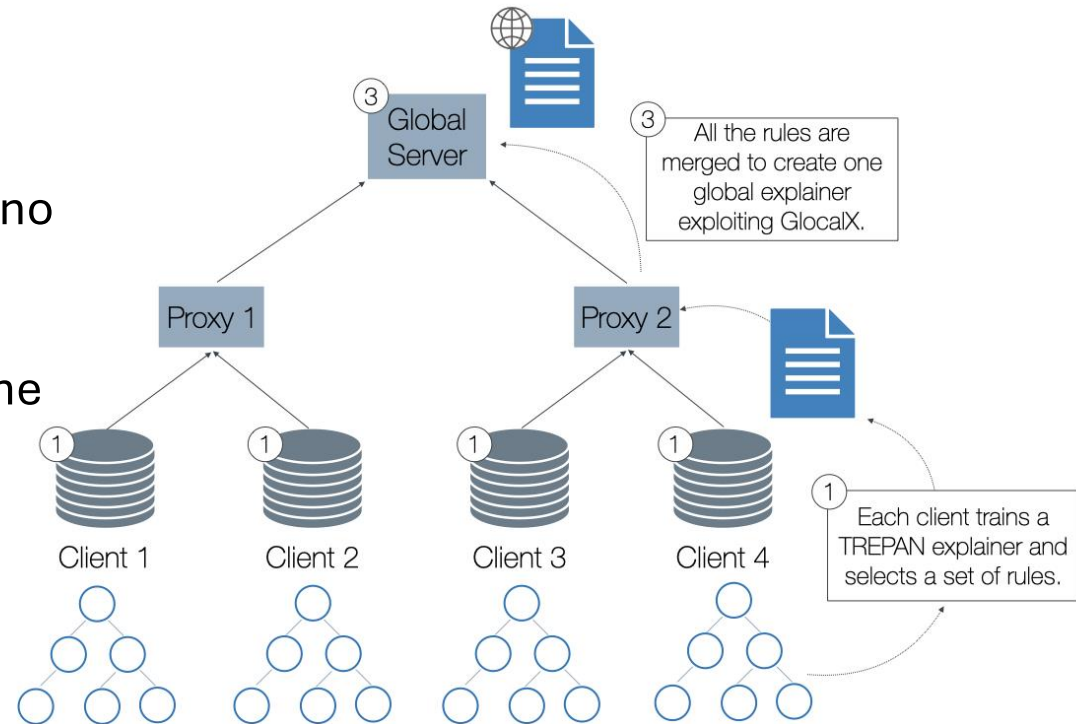
GLOBAL to loCAL eXplainer (**GlocalX**)

It hierarchically merges local explanations.

Explanations with lower fidelity are filtered out.

# How to merge rules: drawback

- GlocalX needs data to perform its tasks.
- But if we use real data outside the clients, we are no longer respecting privacy
- Solution: generate synthetic data that resemble the original ones



# Conclusion

GLOR-FLEX: A local-to-global post-hoc explanation method which generates rules for Federated Learning approaches.

- It uses TREPAN to generate global explanations at the client side;
- It uses GlocalX to merge the rules.

WRAP-UP:

- No private data exploited in the procedure
- Interpretability enhanced

