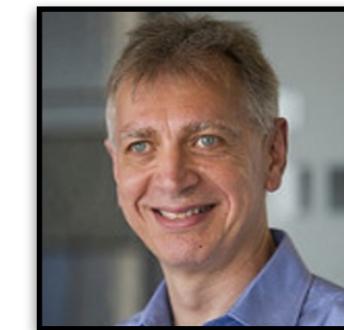


# Privacy-Preserving Federated Data Sharing



Ferdinando Fioretto  
Georgia Tech



Pascal Van Hentenryck  
Georgia Tech

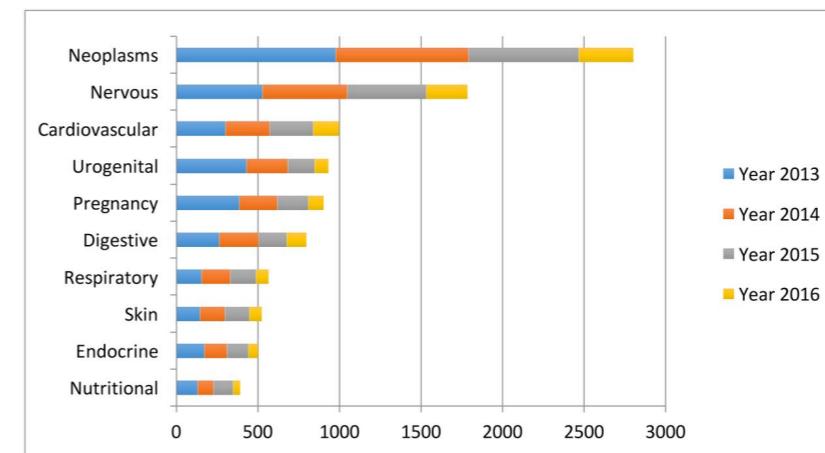
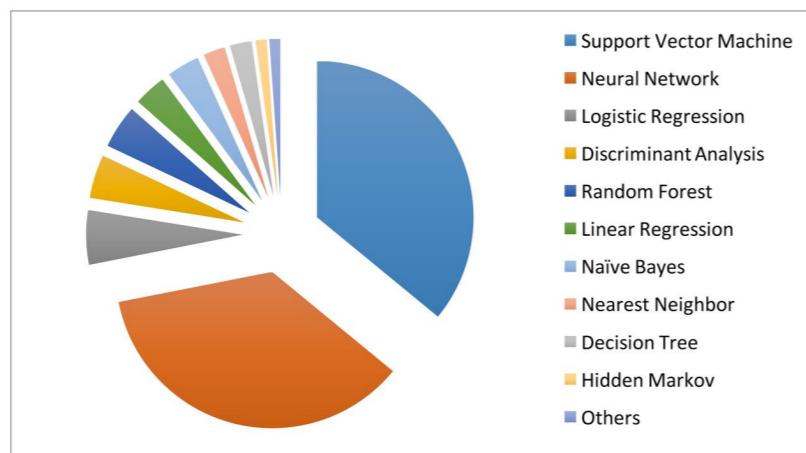
May 15, 2019

# We live in a data-driven world



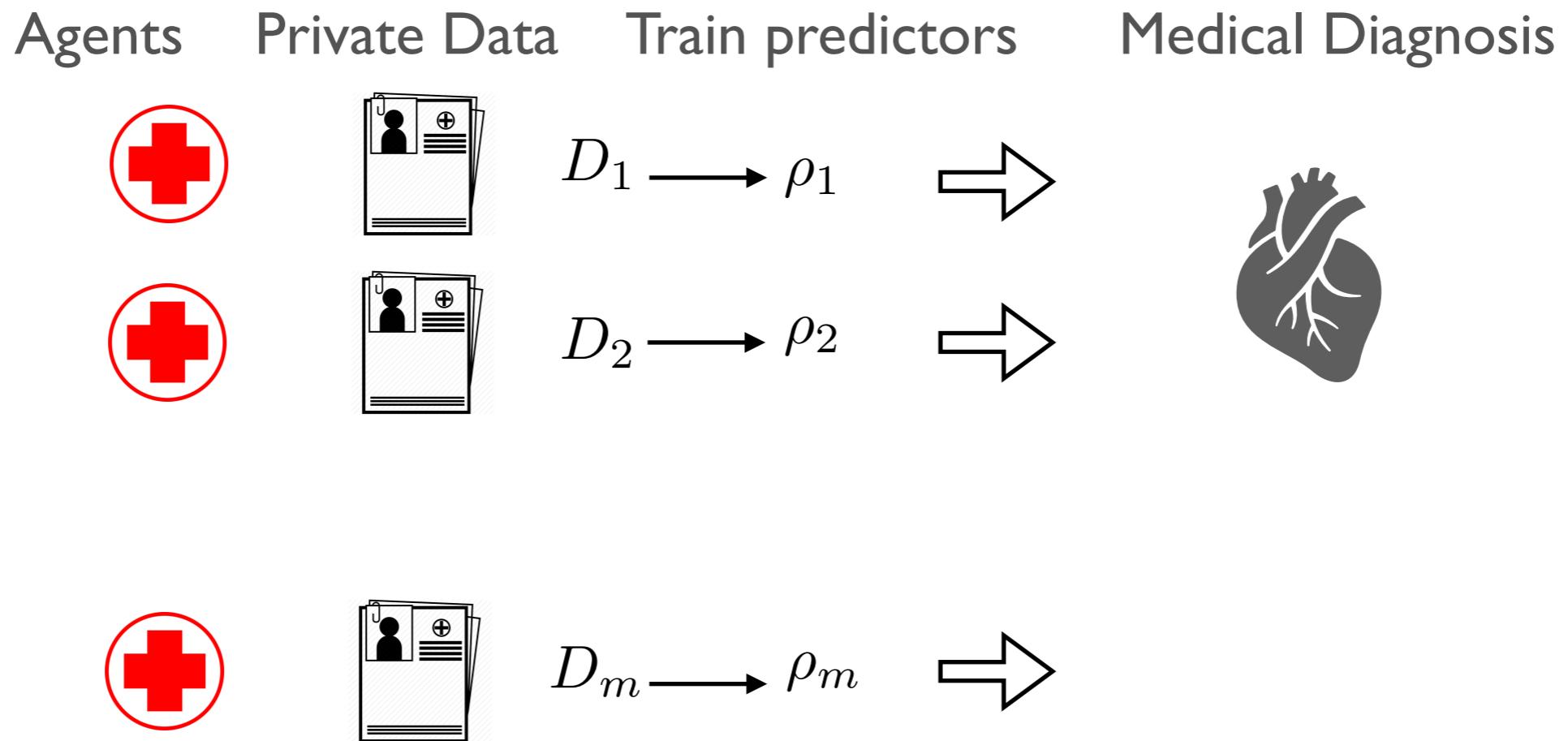
# Motivations

- The availability of personal data is responsible for rapid advancements of AI and ML applications.
- Many data-driven ML models to support e.g., screening, diagnosis, and treatment assignment in healthcare
- Transition from proprietary data acquisition and processing to distributed data ecosystems.
- Personal data is **very sensitive** in many domains (e.g., healthcare)

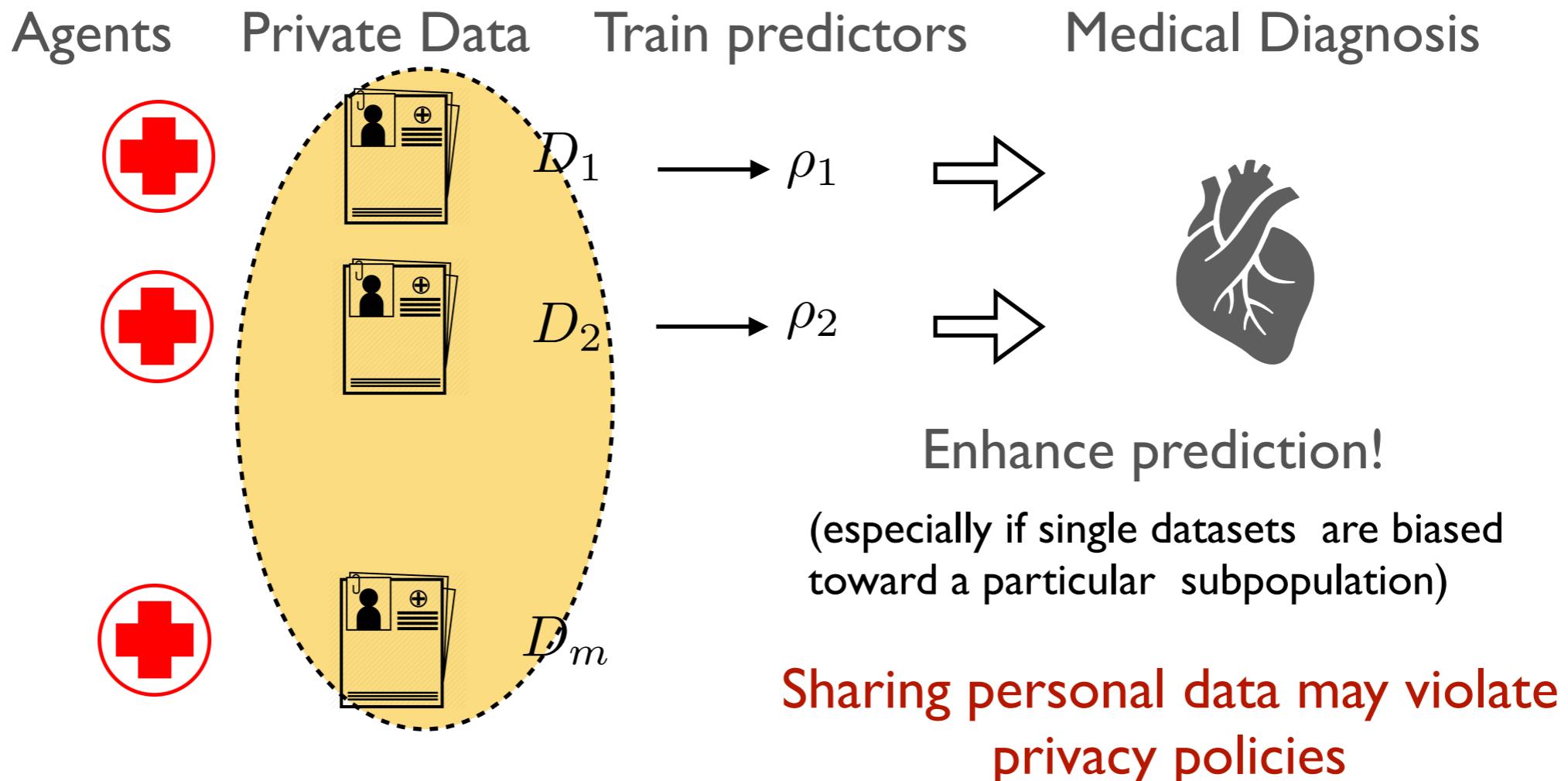


Top ML models adopted (left) and leading disease types (right)  
considered in the ML & Healthcare literature [1]

# The Problem



# The Problem

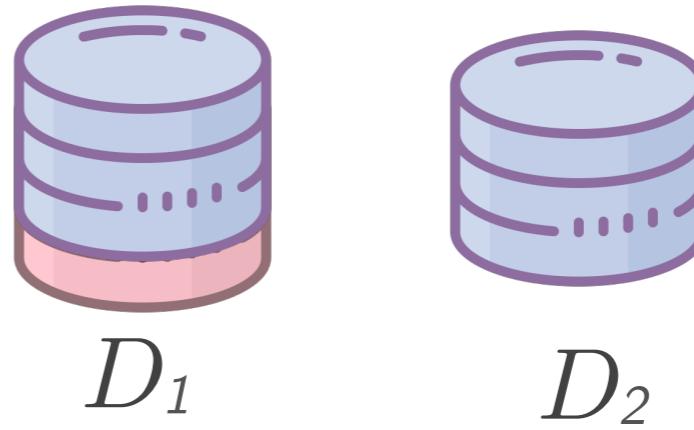


How to build a Federated Data Sharing protocol  
that preserves the individuals data Privacy  
and it is suitable for prediction tasks

# Differential Privacy

For every pair of inputs that  
differs in one row

For every output  $O$



A randomized algorithm  $\mathcal{A}$  is  $\epsilon$ -differentially private if:

$$\frac{\Pr[\mathcal{A}(D_1) = O]}{\Pr[\mathcal{A}(D_2) = O]} \leq \exp(\epsilon)$$

Intuition: adversary should not be able to use output  $O$   
to distinguish between any  $D_1$  and  $D_2$

# Differential Privacy

## Sensitivity Method & Privacy Properties

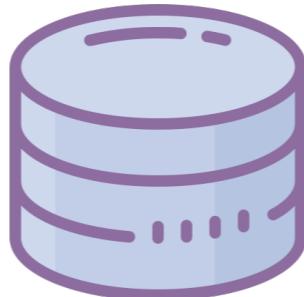
- Consider a query  $Q$  on the dataset  $D$ . The *sensitivity method* injects noise to the output  $Q(D)$  that depends on the **sensitivity of the query**:

$$\Delta_Q = \max_{D_1 \sim D_2} \|Q(D_1) - Q(D_2)\|_1$$

- Laplace Mechanism [2], Exponential Mechanism [3], ...
- **Differential Privacy Important Properties:**
  - **Composability:** Graceful privacy loss degradation when composing multiple DP algorithms.
  - **Immunity to post-processing:** Privacy guarantees are preserved by arbitrary, data-independent, post-processing steps.

# Empirical Risk Minimization

- Labeled data  $D = \{\mathbf{x}_i, y_i\}_{i \in [n]}$



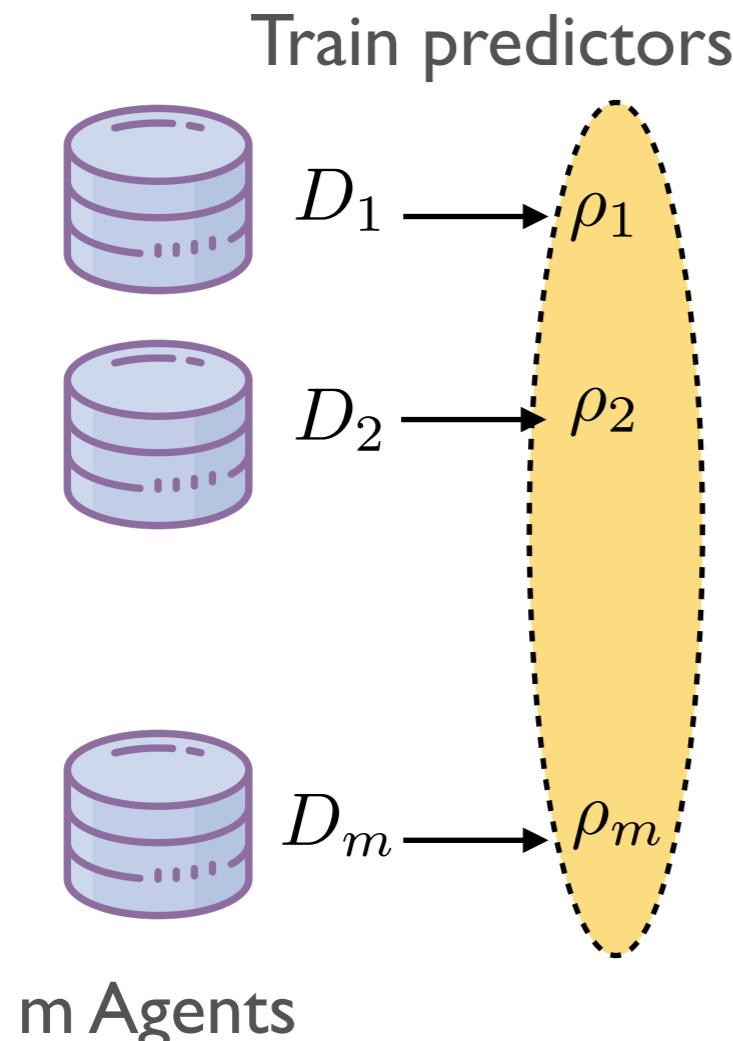
$\underbrace{\mathbf{x}_i \in \mathbb{R}^d}_{(\text{age, gender, ..., diabetes})}$      $y_i \in [-1, 1] \subseteq \mathbb{R}$

- ERM aims at choosing a predictor  $\rho_{\mathbf{w}}$  that minimizes the (regularized) empirical loss:

$$J(\rho_{\mathbf{w}}, D) = \frac{1}{n} \sum_{i=1}^n (\ell(\rho_{\mathbf{w}}(\mathbf{x}_i), y_i)) + \lambda c(\mathbf{w})$$

- Different loss functions allow us to express different predictors. E.g.:
  - Logistic Regression
  - Linear Regression
  - SVM

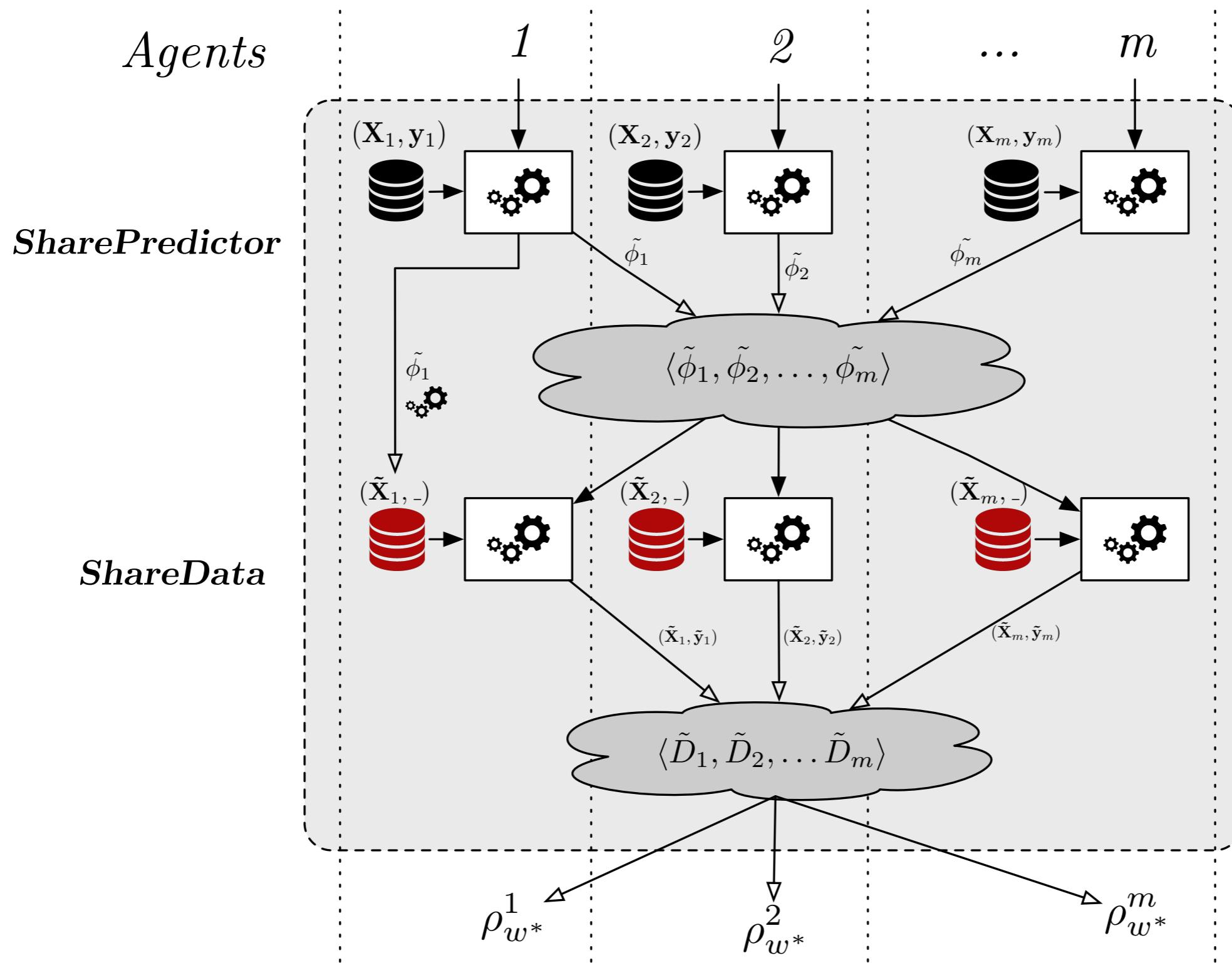
# Multi-agent ERM



- Aggregation schemes:
- [Data] **Majority voting**: Outputs the class that has been predicted more often
- [Predictors] **Committee machine**: Averages the results of each predictor.
- Dataset contain private information

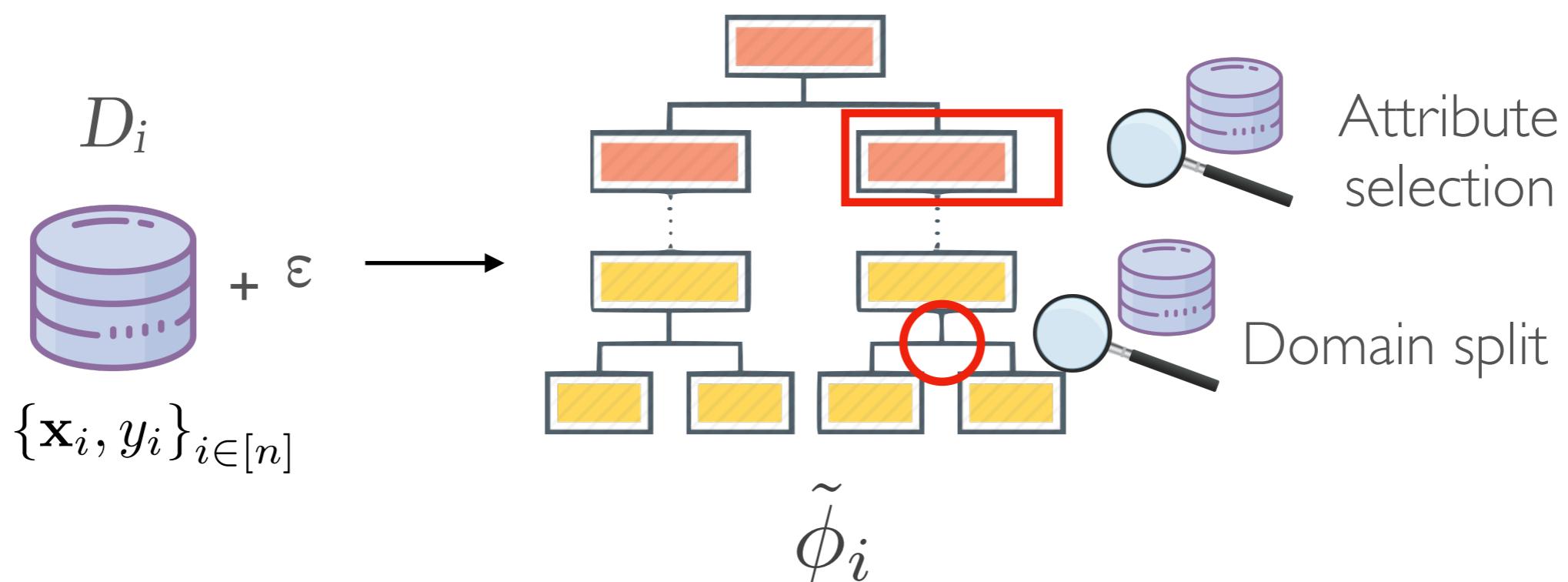
Sharing Data/Predictors  
must be **privacy-preserving!**

# The PFDS Framework

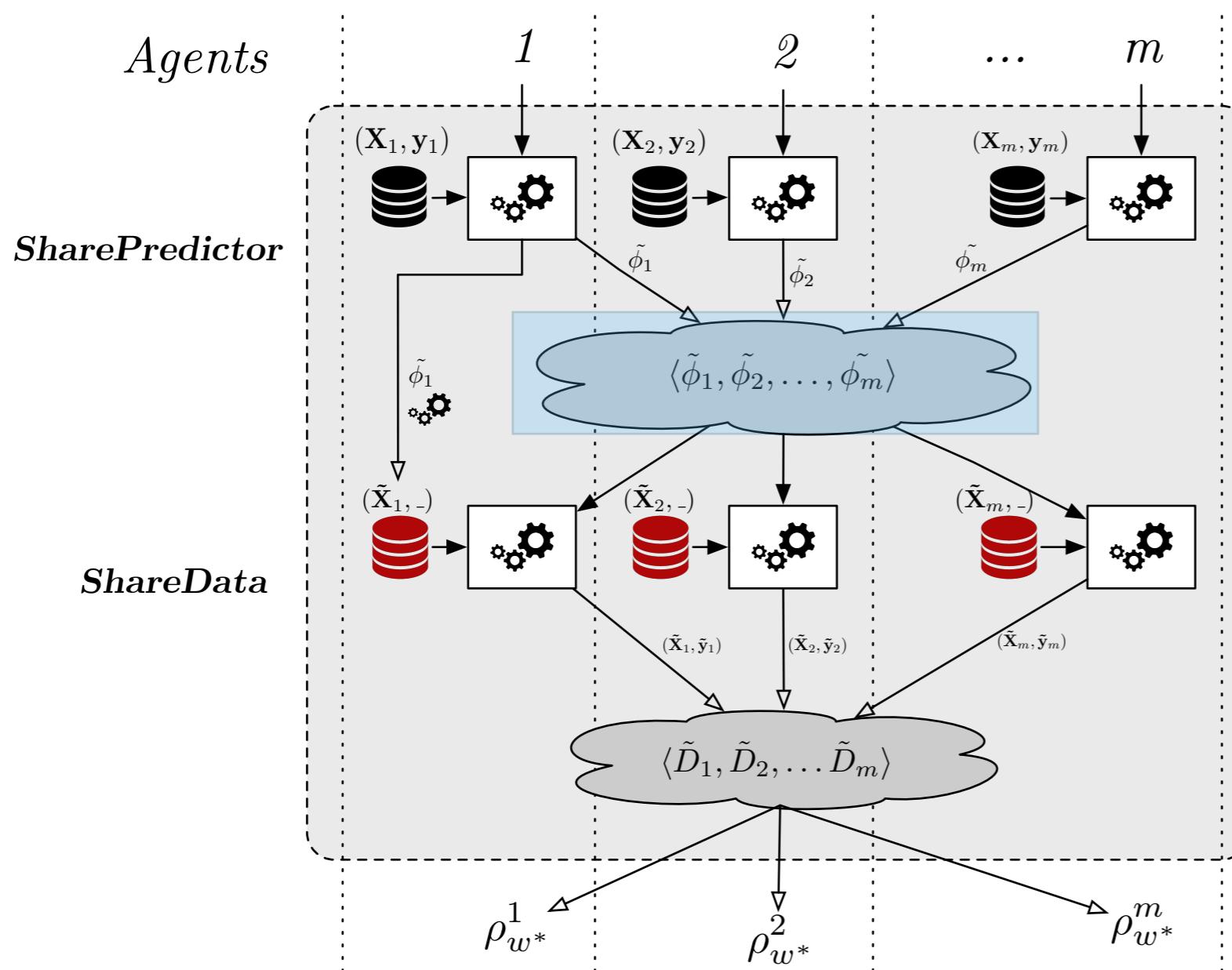


# Predictor Sharing Phase

- Each agent **locally** builds a  $\epsilon$ -differentially private decision tree  $\tilde{\phi}_i$  from its own dataset  $D_i$
- Attribute selection and domain split “look at the data”: require private queries
- Uses a noisy splitting function that relies on sub-sampling the domain of the select value and the Exponential Mechanism



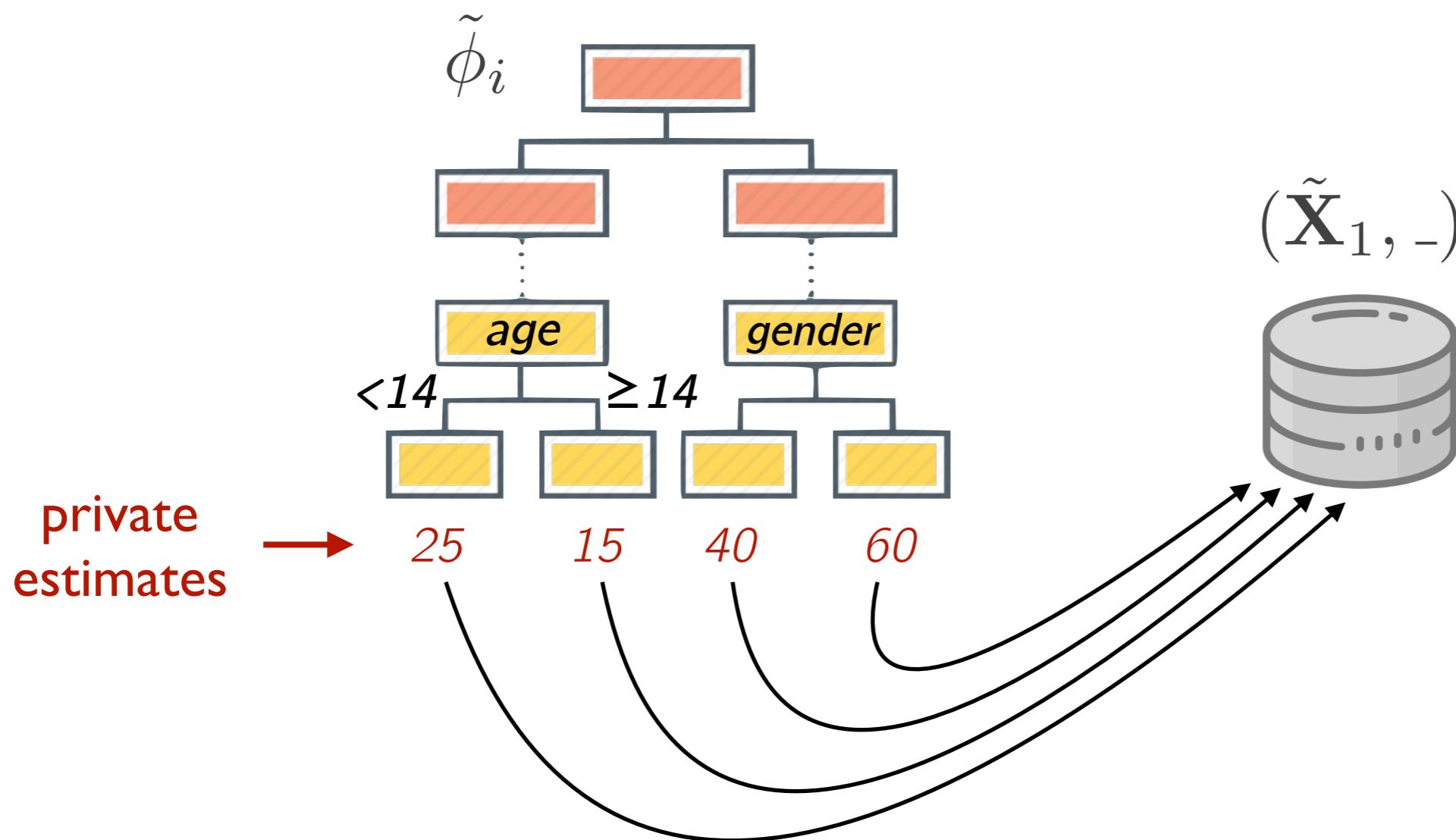
# Predictor Sharing Phase



- The agents share their privacy-preserving predictors

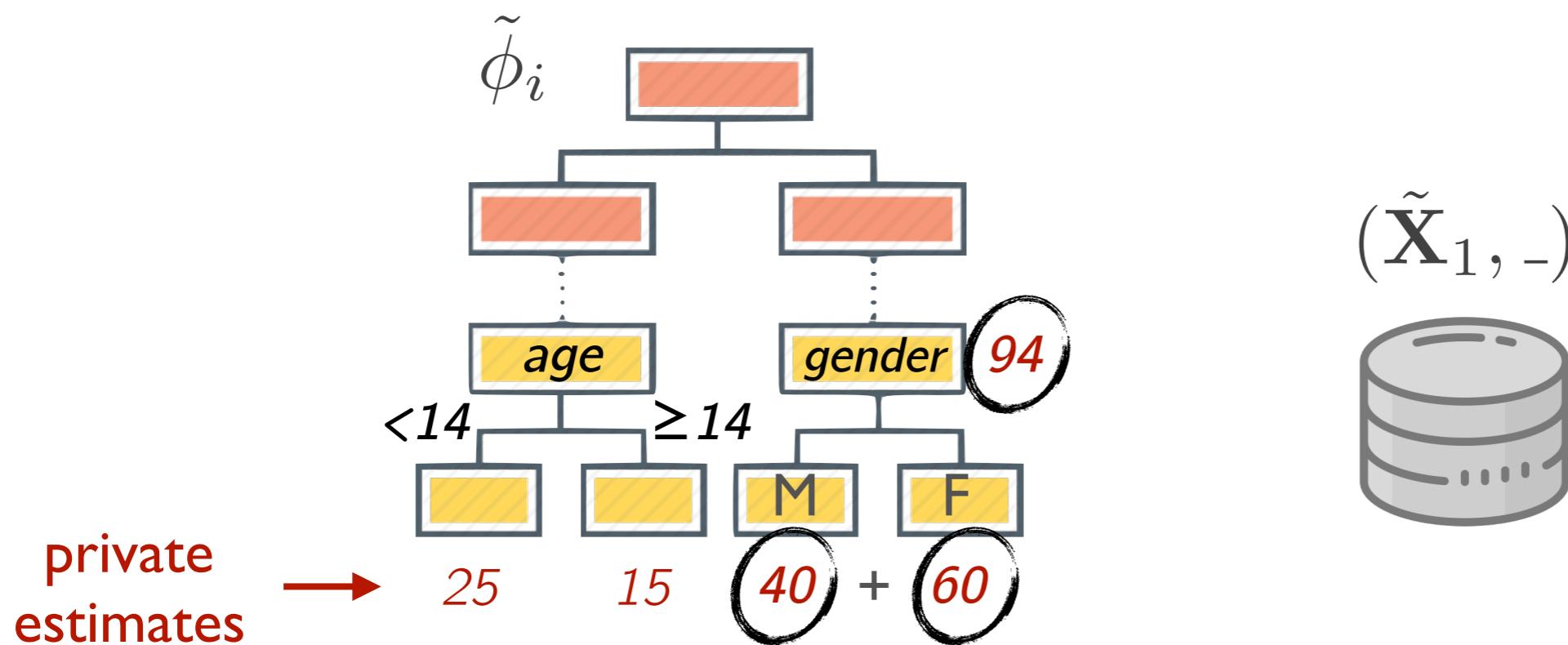
# Data Sharing Phase

- The private predictor  $\tilde{\phi}_i$  is used by the agent  $i$  to create a **unlabeled** version of its dataset.
- The new dataset is generated using only **private estimates** (counting the number of individuals satisfying a given property in the  $X_1$ )



# Data Sharing Phase

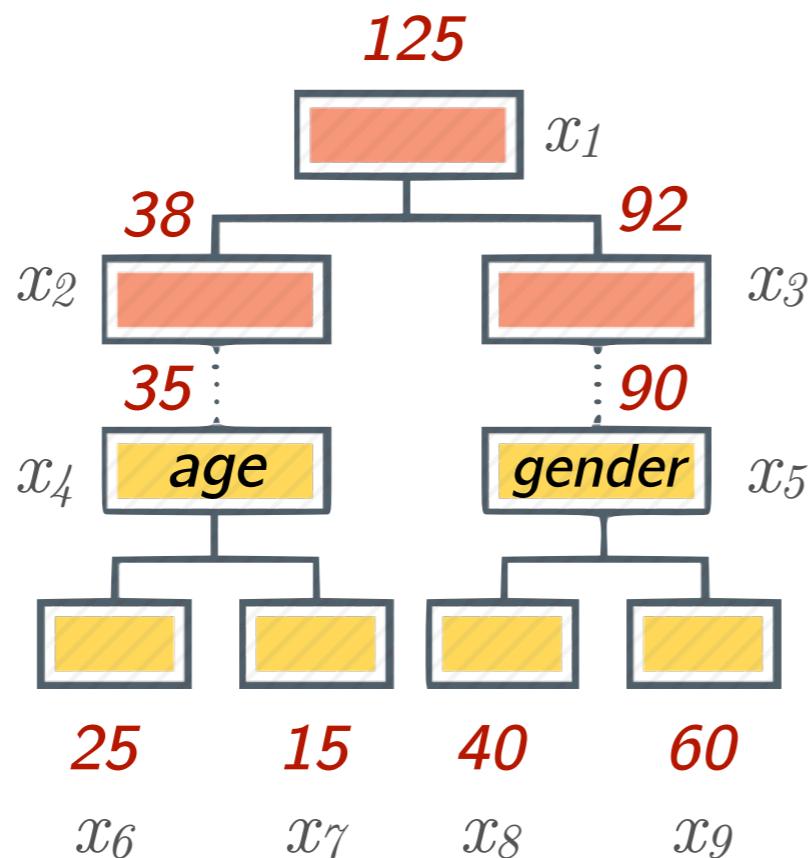
- Because noise is added independently at each node of the tree, hierarchical properties may be inconsistent



# Data Sharing Phase

## Optimization

- To restore consistency, we use convex optimization.
- It redistributed the noise to ensure count consistency.
- THM: The Optimization-based process is  $\epsilon$ -differentially private



$$\min \sum (x_i - \tilde{x}_i)^2$$

$$s.t.: \quad x_i$$

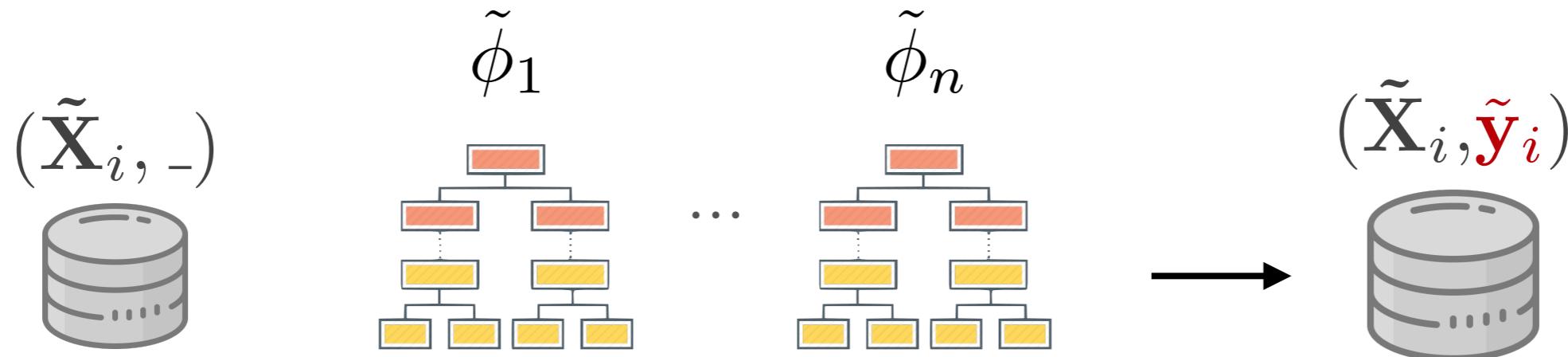
$$\sum_{\tau_i \in ch(\tau_j)} x_i = x_j \quad \forall j \in [9]$$

$$x_i \geq 0 \quad \forall i \in [9]$$

# Data Sharing Phase

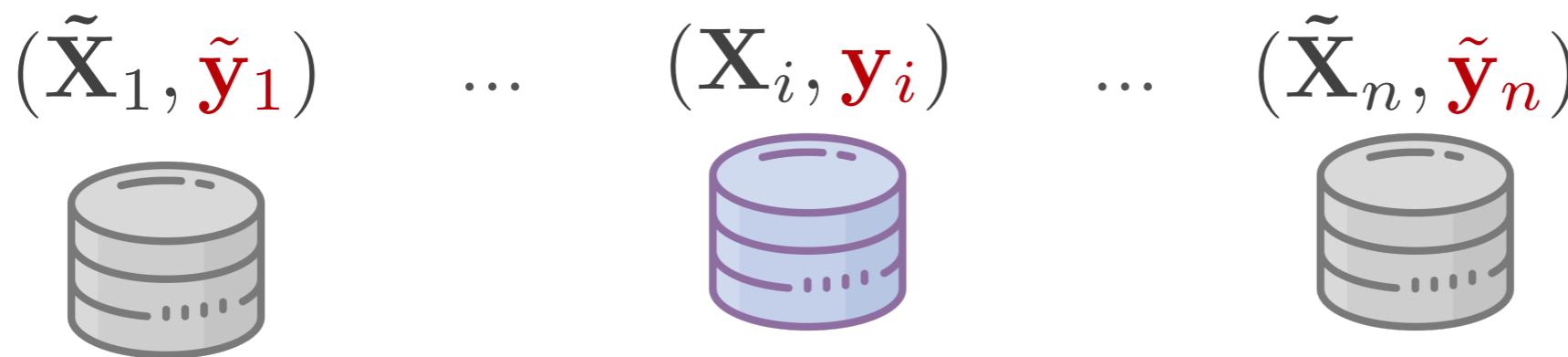
## Labeling Private Data

- Each agent  $i$  uses all shared predictors to find the *labels* for its private dataset.
- Each dataset item  $\tilde{\mathbf{x}}_i$  is labelled with
$$\tilde{\mathbf{y}}_i = \text{majority\_vote}(\tilde{\phi}_1(\tilde{\mathbf{x}}_i) \dots \tilde{\phi}_n(\tilde{\mathbf{x}}_i))$$
- The agent share its privacy-preserving (labeled) data with all other agents



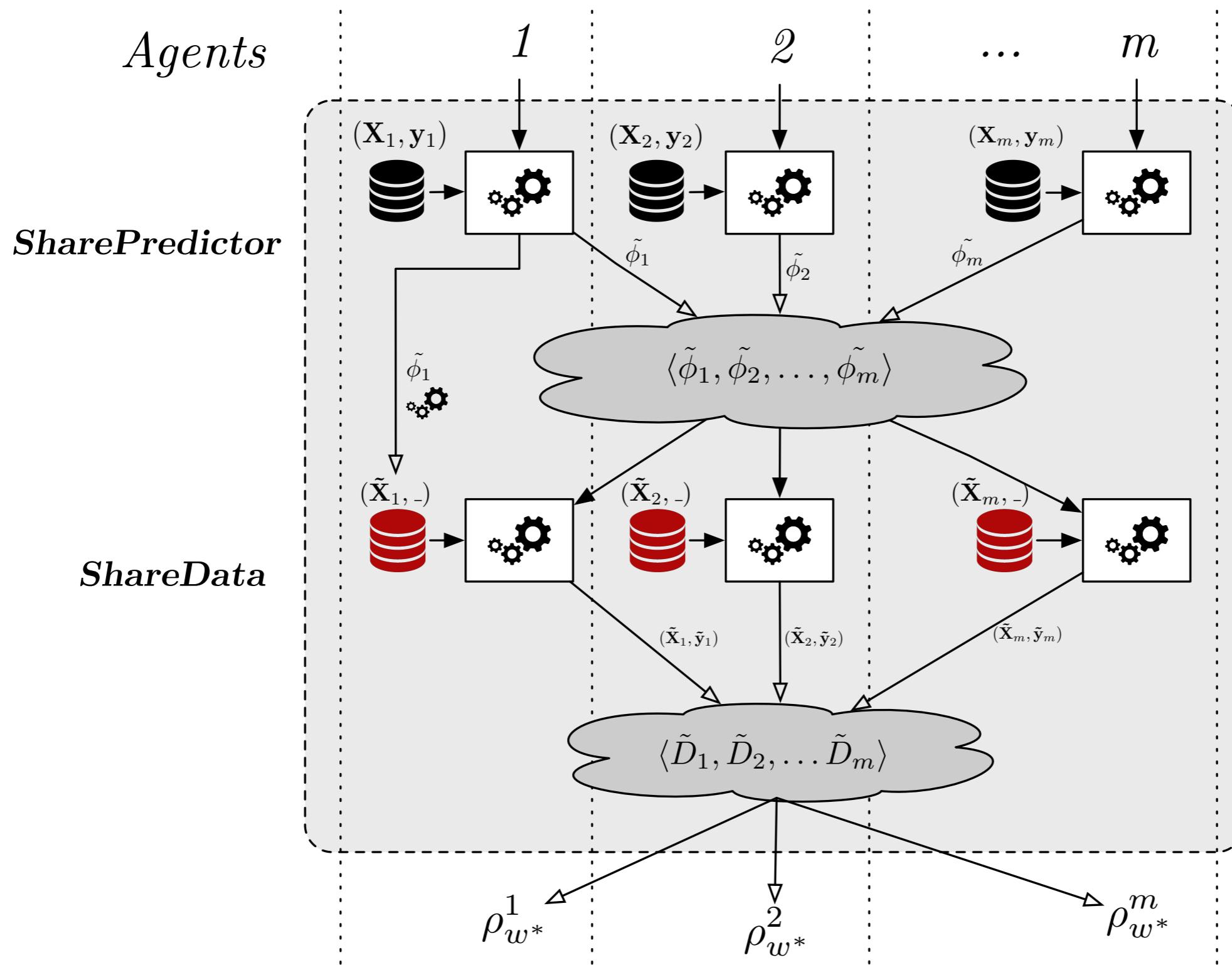
# Training on Federated Data

- Each agent uses all shared privacy-preserving dataset and its own dataset to train its predictor  $\rho_{w^*}^i$



$$J(\rho_{\mathbf{w}}, D) = \frac{1}{n} \sum_{i=1}^n (\ell(\rho_{\mathbf{w}}(\mathbf{x}_i), y_i)) + \lambda c(\mathbf{w}) \quad \text{ERM training}$$

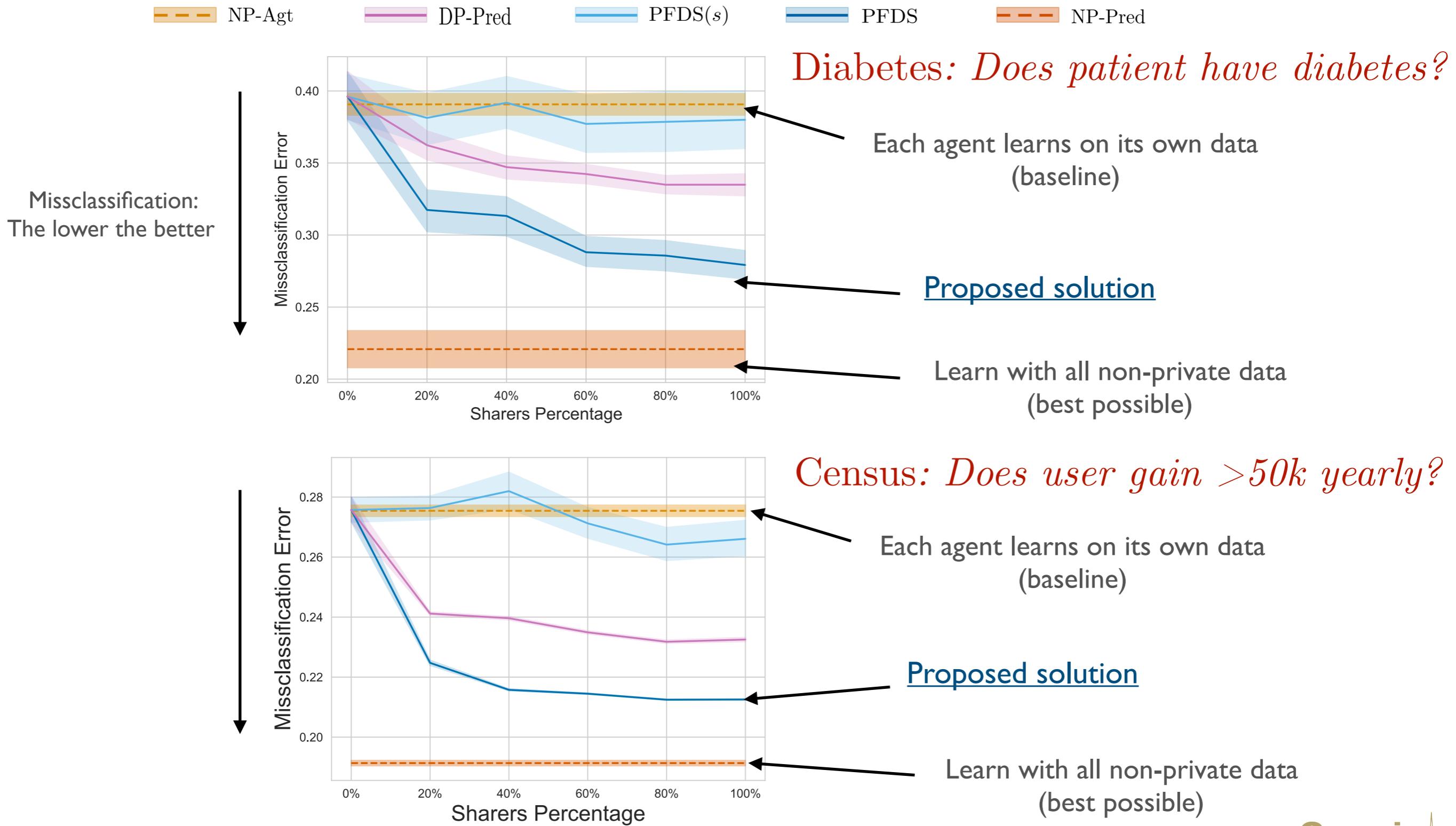
# The PFDS Framework



# Evaluation: Methods

- **PFDS**: method developed in this work
- **PFDS(s)**: PFDS w/out Predictor Sharing (to label data)
- **DP-Pred**: Objective perturbation<sup>[4]</sup> (log. reg. and SVM) or functional mechanism<sup>[29]</sup> (lin. reg.) w/ majority voting
- **NP-Agt**: Non-private predictor executed on the agent data
- **NP-Pred**: Non-private predictor executed on all data

# Evaluation (summary)



# Conclusions

- Companies are increasingly leveraging distributed data.
- Typically, more data implies better models, however privacy regulation prevent data owners to share their data
- We investigated whether sharing private version of the agents datasets may improve the agents predictions
- We proposed a Privacy-preserving Federated Data Sharing protocol to allow agents to share their data privately.
- On Multi-agent ERM tasks, our results show that it is possible to both
  - I. Guarantee privacy and
  2. Learn more accurate models than those learned by single-agent systems on non-private data.

# Conclusions

- Companies are increasingly leveraging distributed data.
- Typically, more data implies better models, however privacy regulation prevent data owners to share their data
- We investigated whether sharing private version of the agents datasets may improve the agents predictions
- We proposed a Privacy-preserving Federated Data Sharing protocol to allow agents to share their data privately.
- On Multi-agent ERM tasks, our results show that it is possible to both
  1. Guarantee privacy and
  2. Learn more accurate models than those learned by single-agent systems on non-private data.

Ask us any additional question at the poster session

Thank you!



[fioretto@gatech.edu](mailto:fioretto@gatech.edu)  
[pvh@isye.gatech.edu](mailto:pvh@isye.gatech.edu)

# References

- [1] Jiang, Fei, et al. "Artificial intelligence in healthcare: past, present and future." *Stroke and vascular neurology* 2.4 (2017): 230-243.