

UNIVERSITY OF BOLOGNA
SECOND CYCLE DEGREE IN COMPUTER ENGINEERING
Cybersecurity M Project

Dataset Anonymizer

AI-based system for anonymizing data, ensuring privacy while maintaining reliability for statistical analysis

Authors:

Francesco Simoni

Giacomo Vigarelli

Data:

16/12/2024

ACADEMIC YEAR 2024-2025

Project Goals



I Project Goals

Developing an AI model that ensures the anonymization of a dataset with good accuracy

The following will be presented:

- Configuring, training, and evaluating ML models on the MNIST dataset
- Integrating Differential Privacy mechanisms
- Analyzing accuracy and loss metrics for models with different privacy levels



Differential Privacy



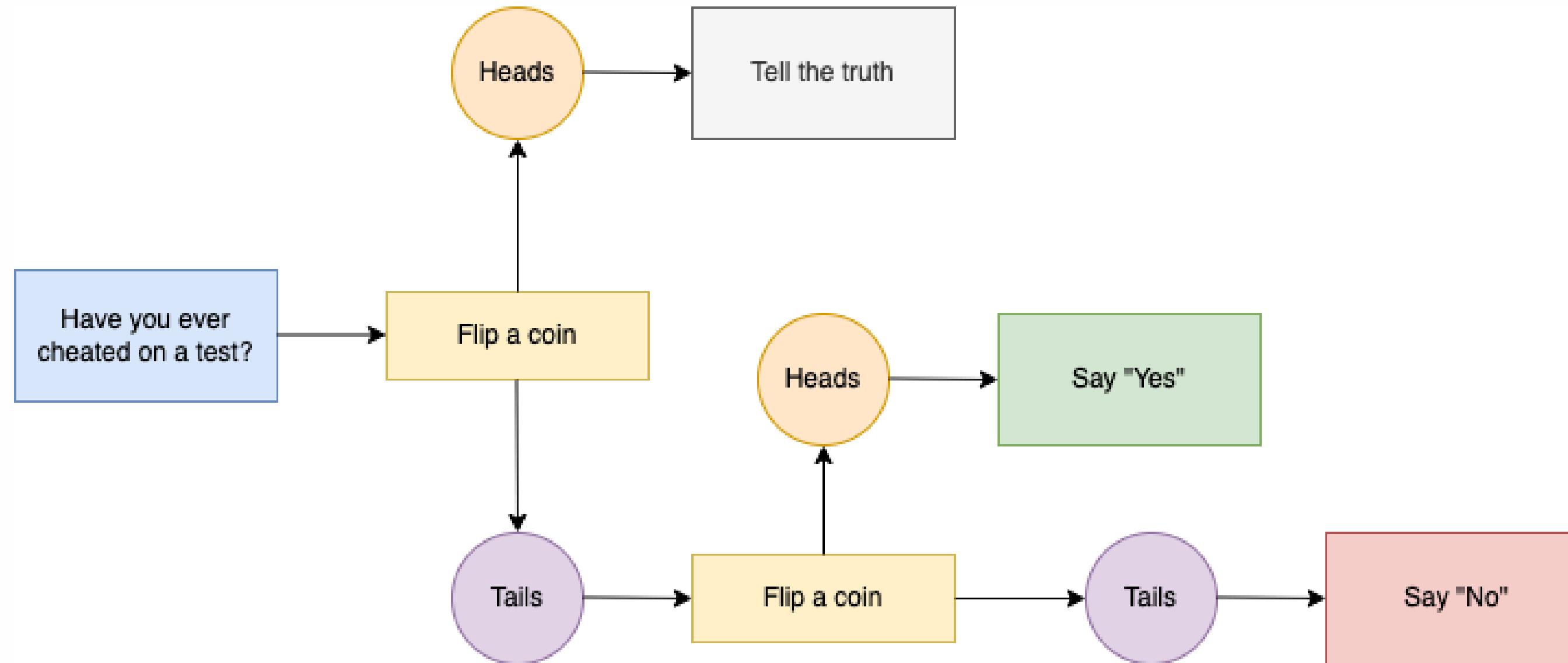
II Differential Privacy

Key Concepts in Differential Privacy

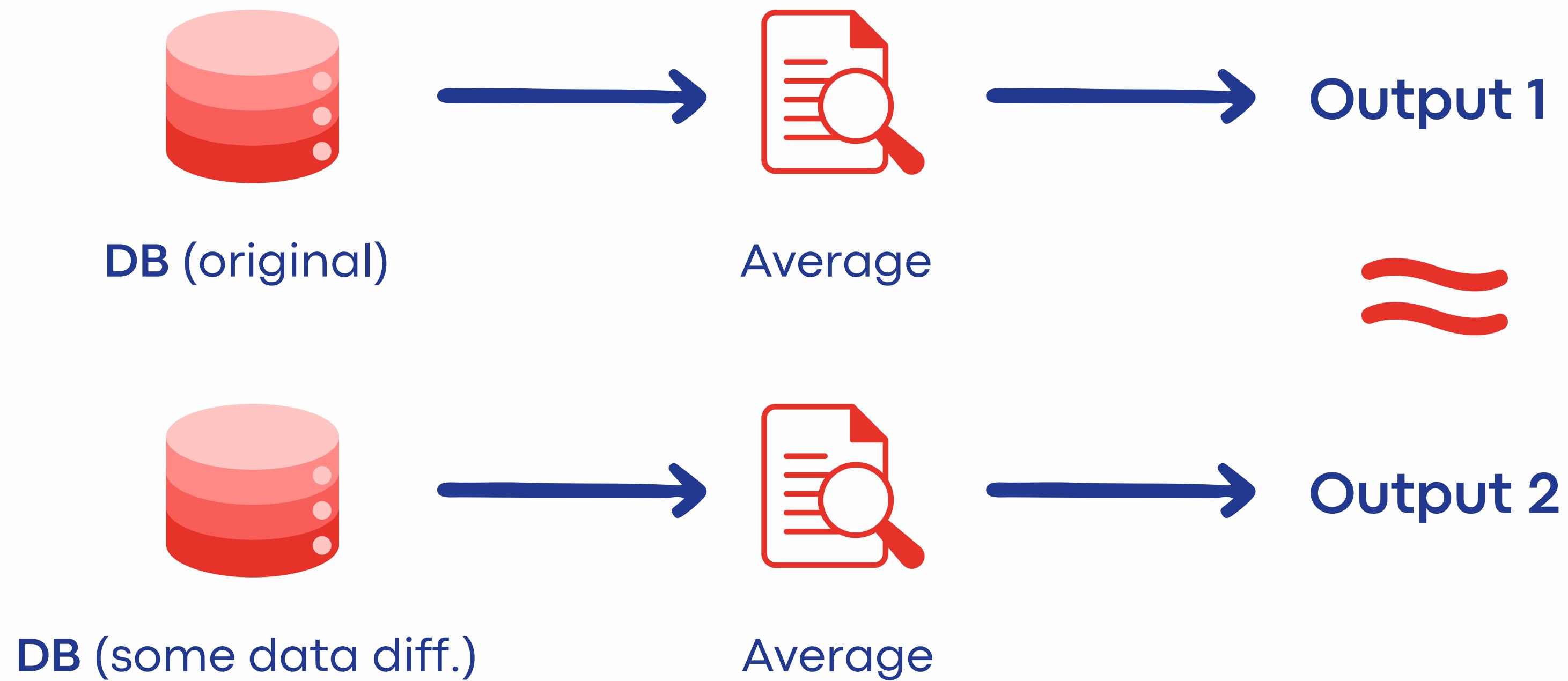
- Technique used to protect Individual Privacy in Dataset
- Ensure that presence/absence of specific data doesn't change the output of a statistic analysis
- If an attack query is done, it is nearly impossible to determine the specific contribution



II Example



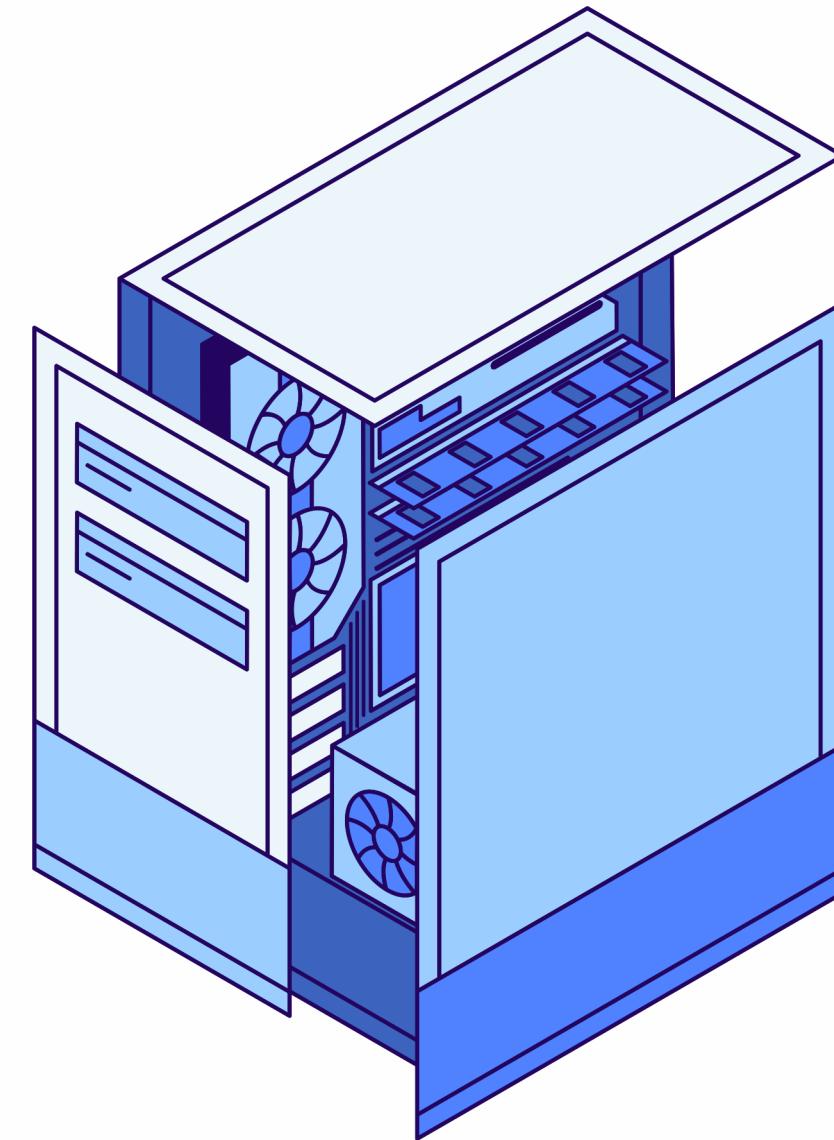
II Example



II Epsilon Parameter

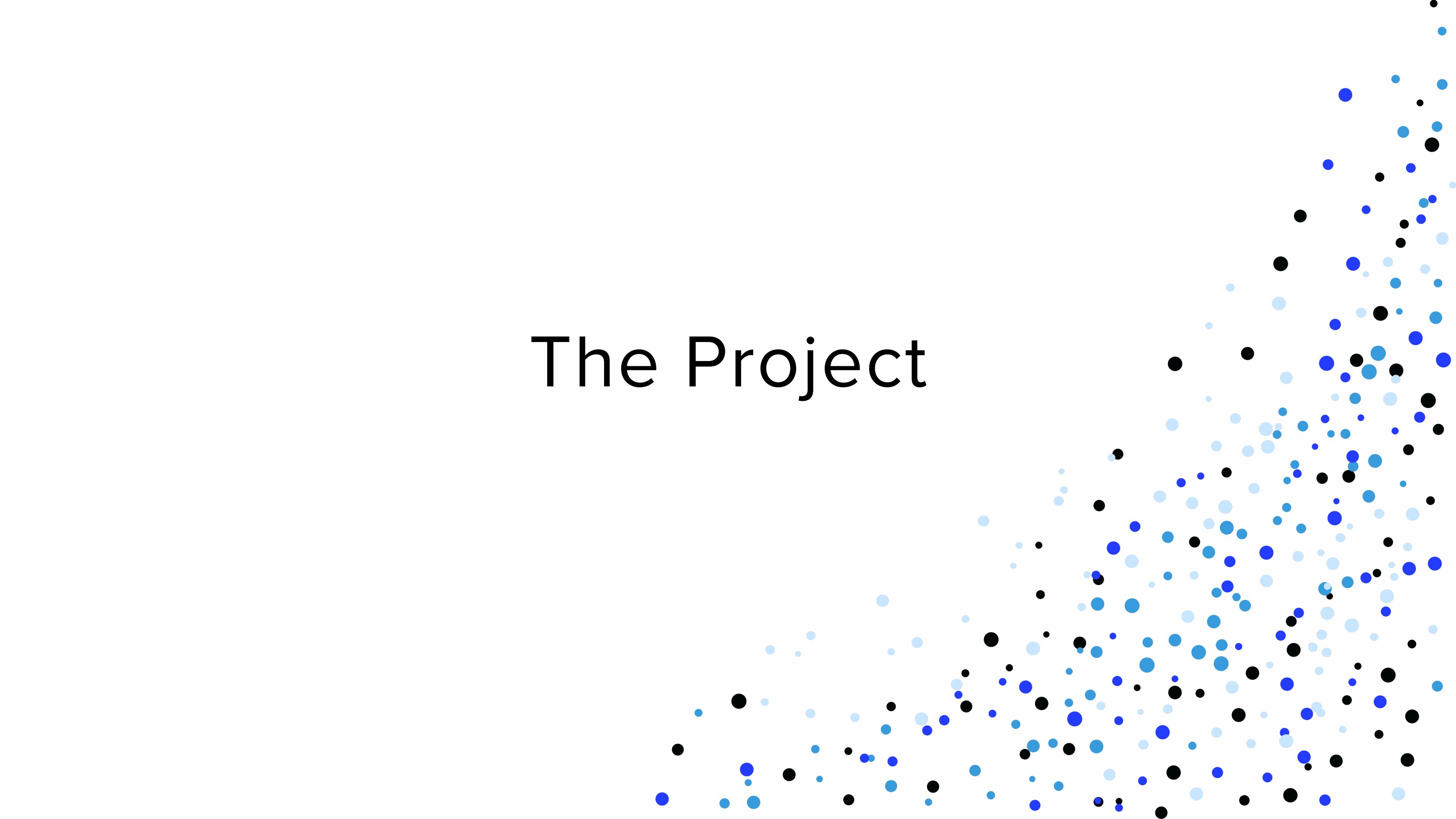
Properties of ϵ , derived from the DP formula:

- Measures privacy level protection
- Low ϵ : Higher privacy (Lower accuracy)
- High ϵ : Higher accuracy (Lower privacy)
- Different ϵ tests to find optimal balance



The parameter will be fundamental for the analysis and final results.

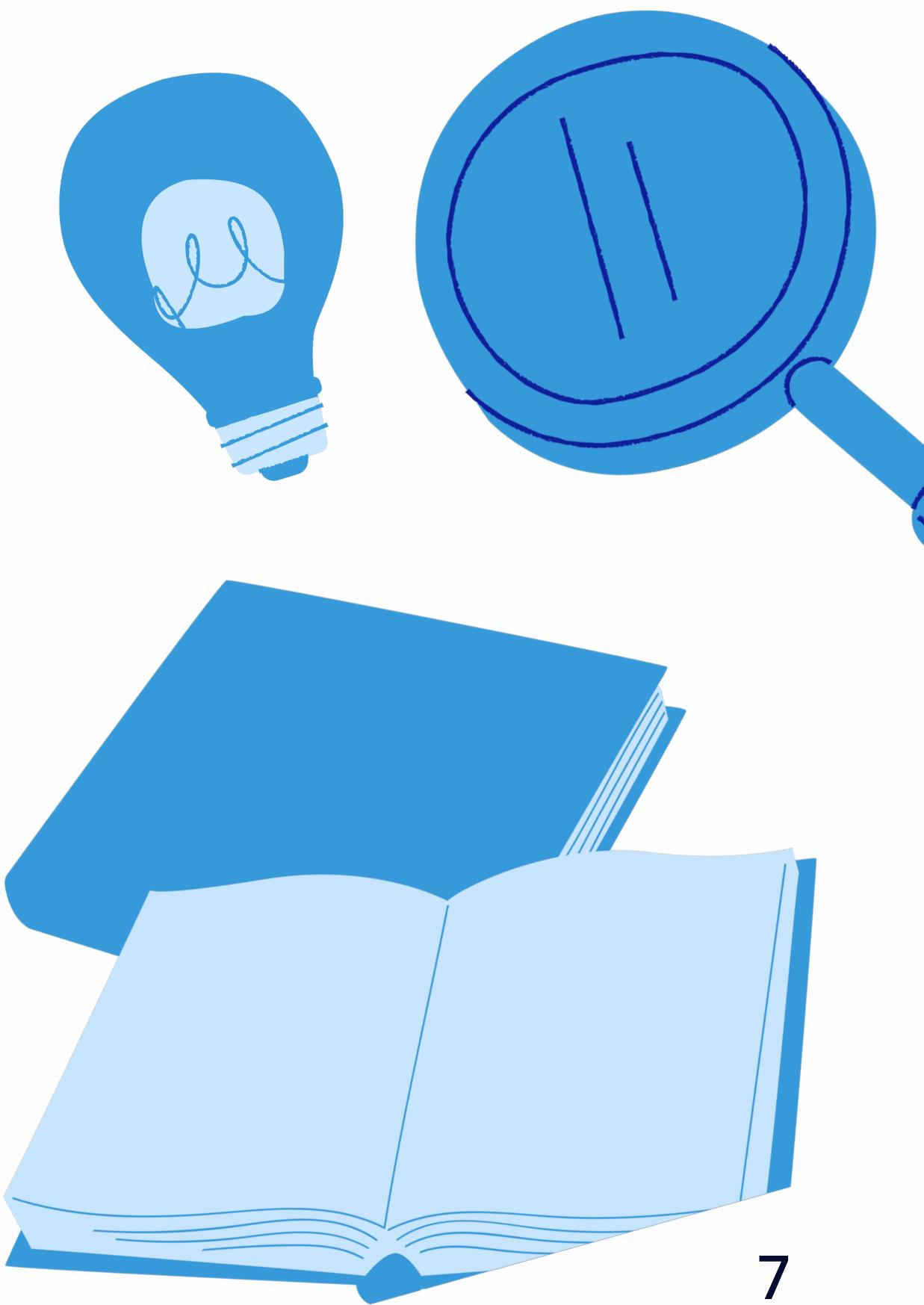
The Project



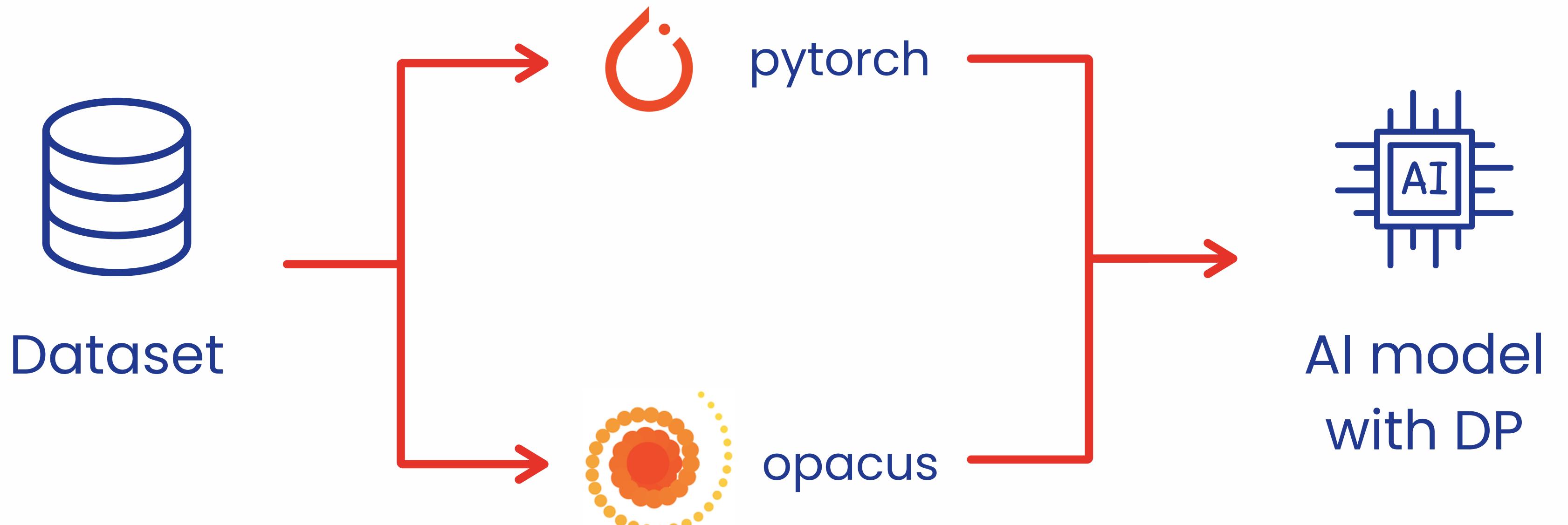
III The Project

Technologies Used

- **Python 3.11**: the most common programming language for ML models
- **PyTorch**: a training framework for Python along with *TorchVision*, a specific extension for the MNIST dataset
- **Opacus**: a library to implement DP in PyTorch-based models
- **Matplotlib**: a library for graph visualization



III Implementation



III Dataset

Which dataset should we use?

We evaluated two datasets:

- **MNIST**: An image dataset containing 60,000 samples of handwritten digits (0 to 9)
- **ADULTS (UCI)** : A textual dataset with 48,000 records, each containing an individual's name, surname, annual gross income and other sensible data

A grid of handwritten digits from the MNIST dataset, showing rows of 0s, 1s, 2s, 3s, 4s, 5s, 6s, 7s, 8s, and 9s.

0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9 9 9

We chose *MNIST* due to technical incompatibilities with textual datasets and the tools used.

III Project Modular Structure

We designed a *modular* structure in Python to ensure that each .py file has a clear and specific responsibility. This approach improves **code maintainability** and **scalability**. The **src** folder contains the following files:

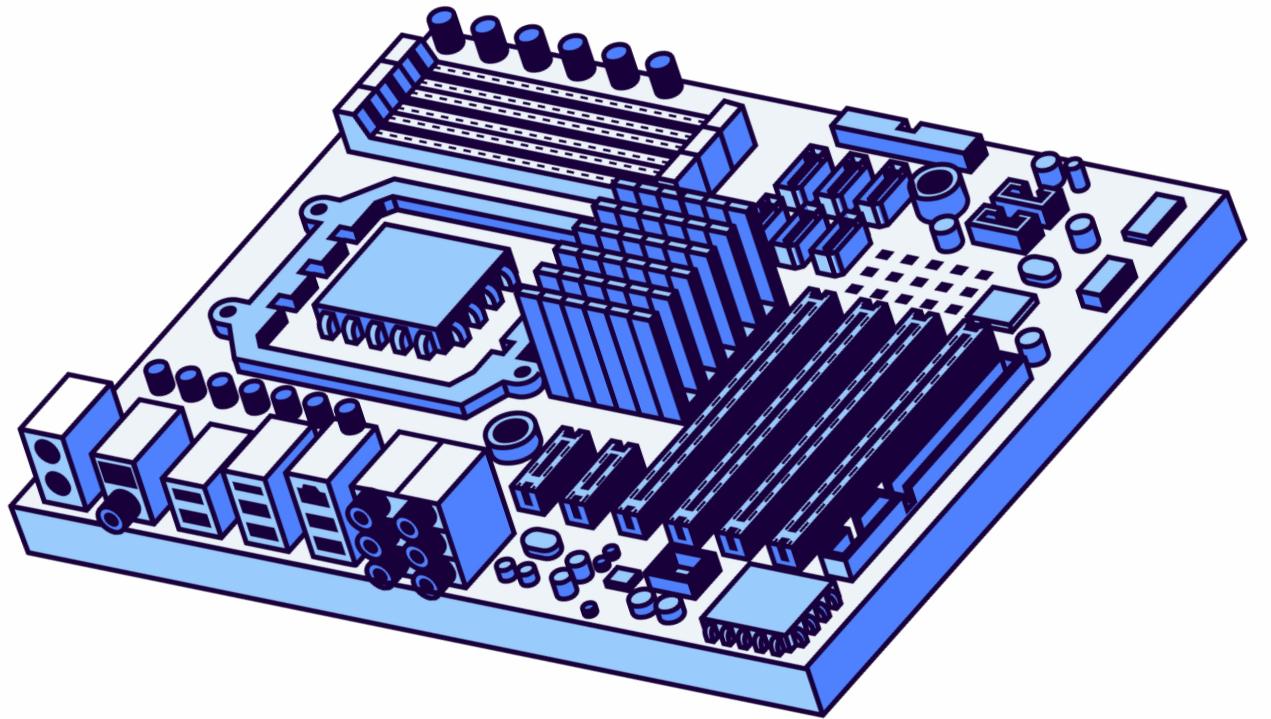
```
src
|
|-- plots                      # The resulting graphs are saved
|-- model.py                   # Define and validate the ResNet model
|-- data.py                    # Load and prepare the MNIST dataset
|-- plotter.py                 # Generate metric visualizations
|-- trainer.py                 # Train and evaluate the ML model
|-- main.py                     # Orchestrate all process
```

III Model

ResNet18: A convolutional neural network (CNN) designed for image classification. Its efficiency and architecture make it suitable for tasks with limited computational resources

Modifications:

- Adapted the first convolutional layer to handle grayscale images (single-channel input)
- Configured for 10 output classes to align with the MNIST dataset



III Trainer

The training pipeline implements the integration of DP using the **Opacus** framework. The PrivacyEngine ensures privacy by adding noise to gradients and clipping them, balancing privacy and accuracy. *Custom parameters* like epsilon and epochs can be insert easily

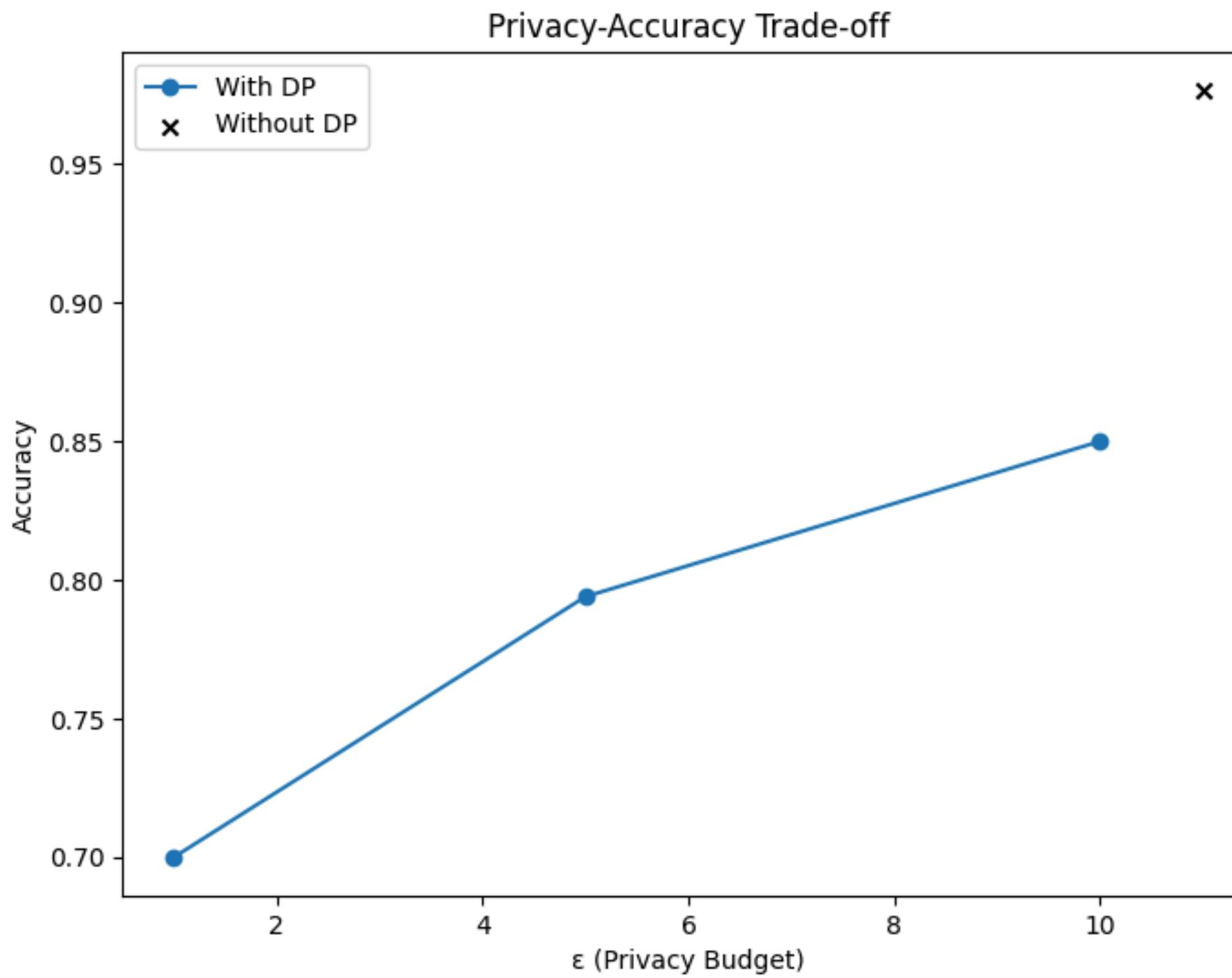
```
1 # Enables DP for the specific model
2 def _enable_dp(self):
3     self.privacy_engine = PrivacyEngine()
4     self.model, self.optimizer, self.train_loader = self.privacy_engine.make_private_with_epsilon(
5         module=self.model,
6         optimizer=self.optimizer,
7         data_loader=self.train_loader,
8         epochs=self.epochs,
9         target_epsilon=self.epsilon,
10        target_delta=self.delta,
11        max_grad_norm=self.clipping_threshold,
12    )
13    print(f"[DEBUG] Using sigma={self.optimizer.noise_multiplier} and C={self.clipping_threshold}")
```

How Opacus integrates with PyTorch.

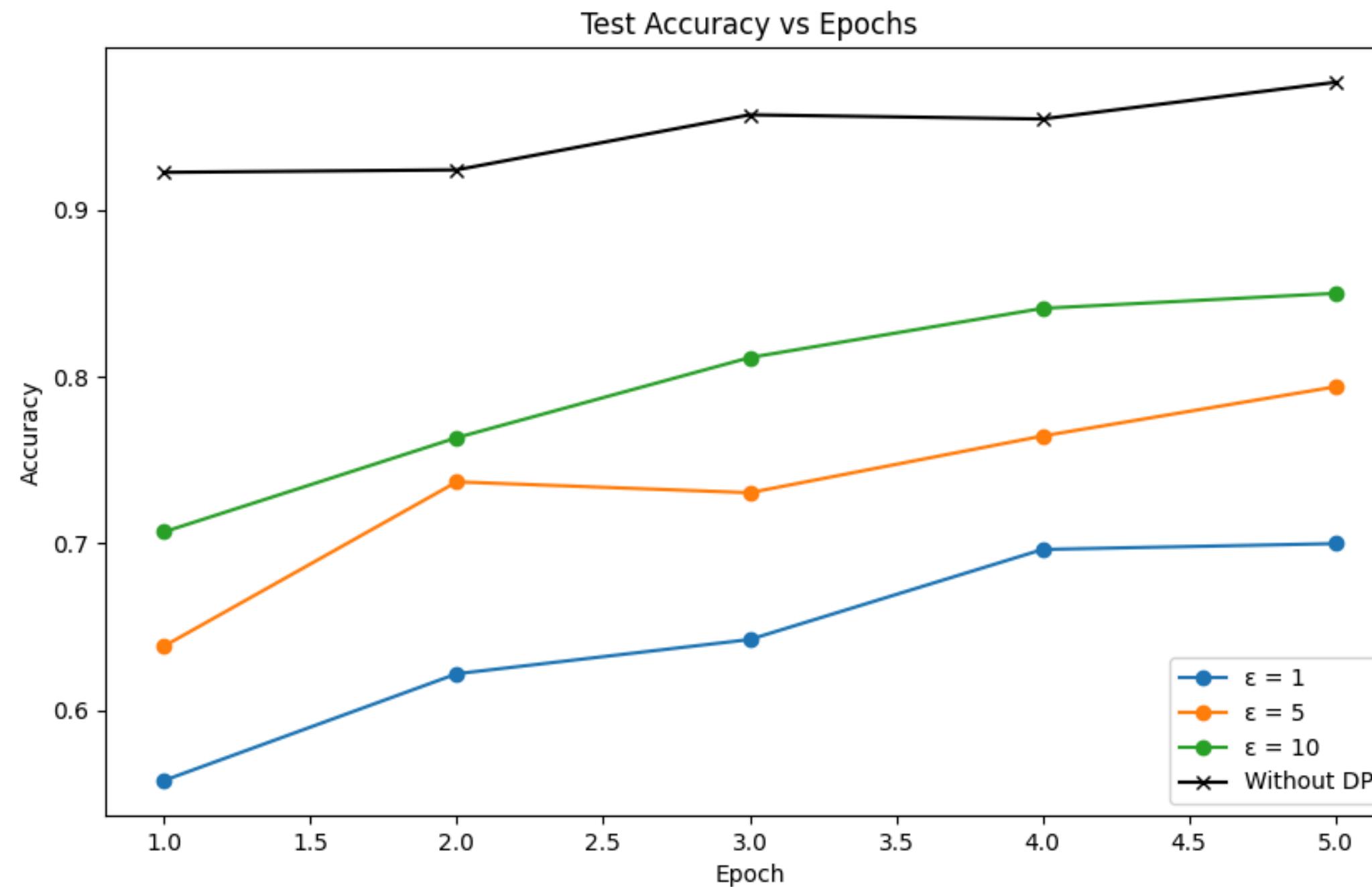
Results & Conclusions



IV Result: *Trade-OFF*



IV Result: Accuracy vs Epochs



V Conclusions

- Successfully implemented a ML prototype for data anonymization using differential privacy
- Training configuration included 5 epochs, [1, 5, 10] epsilon values and reduced datasets MNIST (10,000 training images and 2,000 test images) to optimize computational resources
- There is a clear trade-off between privacy and performance: as privacy weakened (higher ϵ), the model's accuracy improved
- Showcased the importance of customizing existing models to meet specific requirements

16/12/2024

MSc Computer Engineering, Unibo

Thanks for the attention

Authors:

Francesco Simoni

Giacomo Vigarelli