# An Equivalence Between Data Poisoning and Byzantine Gradient Attacks

Sadegh Farhadkhani, Rachid Guerraoui, Lê-Nguyên Hoang and Oscar Villemaud
IC EPFL

EPFL

Global and Personalized Learning

# Global and Personalized Learning

$D_1$

$D_6$

$D_2$

$D_5$

$D_3$

$D_4$

# Global and Personalized Learning

$\theta_1$ $D_1$ $\qquad\qquad\qquad$ $D_6$ $\theta_6$

$\theta_2$ $D_2$ $\qquad\qquad$ $\rho$ $\qquad\qquad$ $D_5$ $\theta_5$

$\theta_3$ $D_3$ $\qquad\qquad\qquad$ $D_4$ $\theta_4$

# Global and Personalized Learning

$\theta_1$ $D_1$

$D_6$ $\theta_6$

$\rho$

$\theta_2$ $D_2$

$$\text{Loss} = \sum_n \mathcal{L}_n(\theta_n, \mathcal{D}_n) + \sum_n \mathcal{R}(\rho, \theta_n)$$
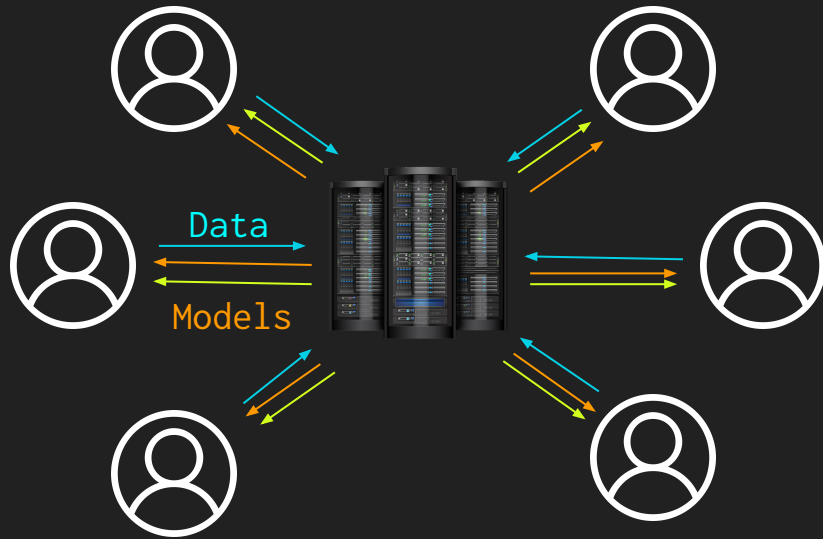
$D_5$ $\theta_5$
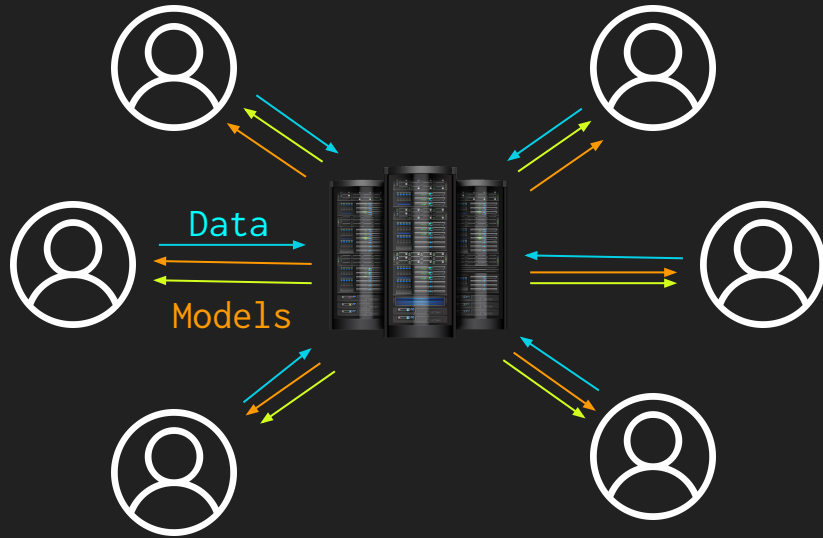
$\theta_3$ $D_3$

$D_4$ $\theta_4$

# Two computation models

## Central computing



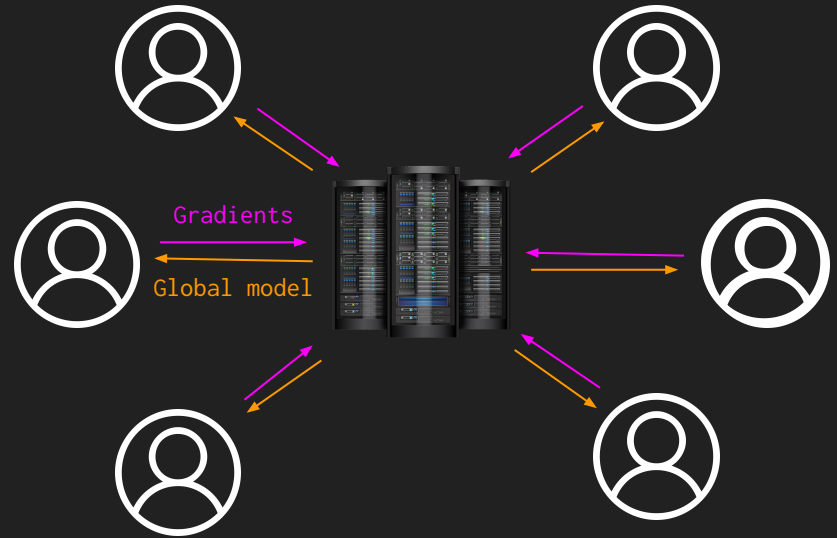Users send data, and receive global and personalized models.

# Two computation models

## Central computing



Data

Models

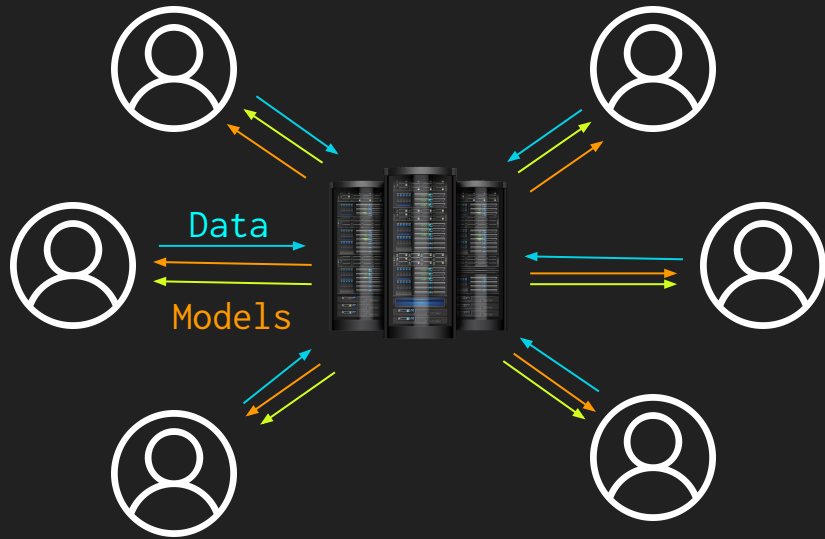Users send data, and receive global and personalized models.

## Federated learning

Gradients

Global model

Users send gradients, and receive global models. They compute personalized models themselves.
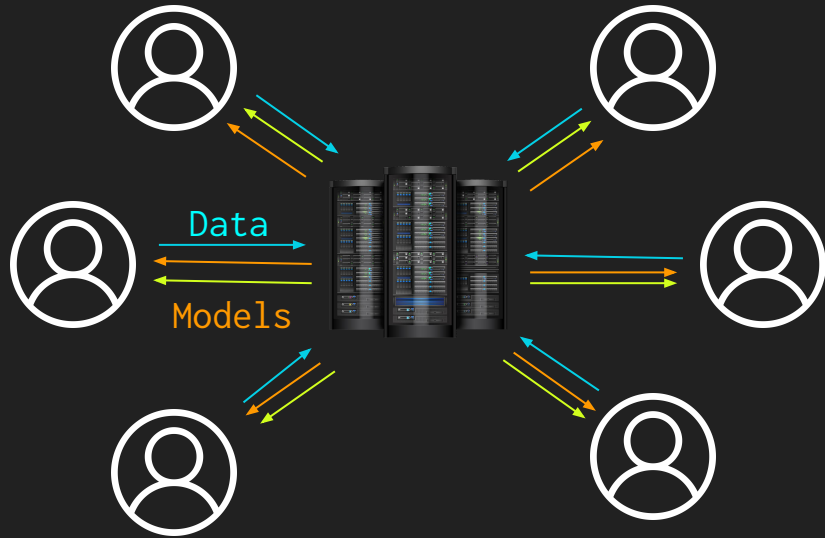
# Two computation models

## Central computing



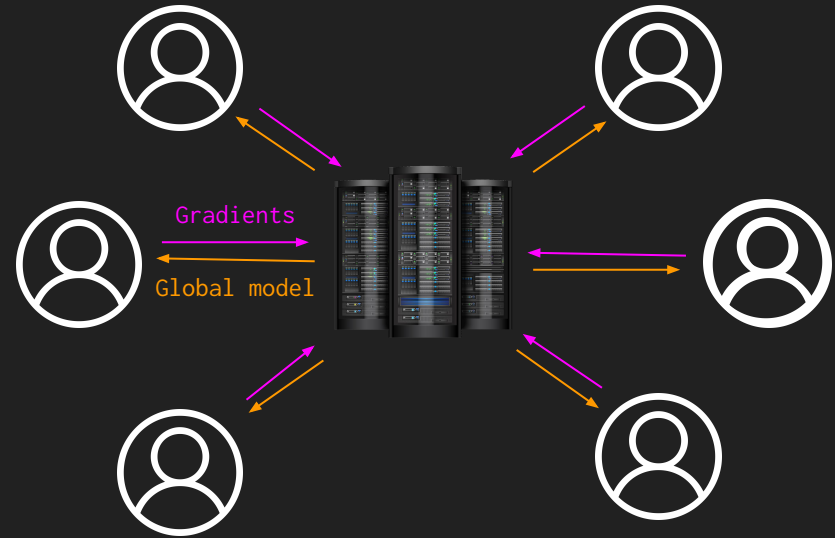Data

Models

Users may send poisonous data.

## Federated learning



Gradients

Global model

# Two computation models

## Central computing



Data

Models

Users may send poisonous data.

## Federated learning



Gradients

Global model

Users may send Byzantine gradients.

# Main theorem

Under realistic and desirable GPL assumptions + convexity, **data poisoning** and **Byzantine gradients** are **equivalent**.

# Our results are very practical!

We develop a new **targeted gradient attack** which successfully **relabels all data**, and we turn it into **effective data poisoning** with **surprisingly few injected data**.
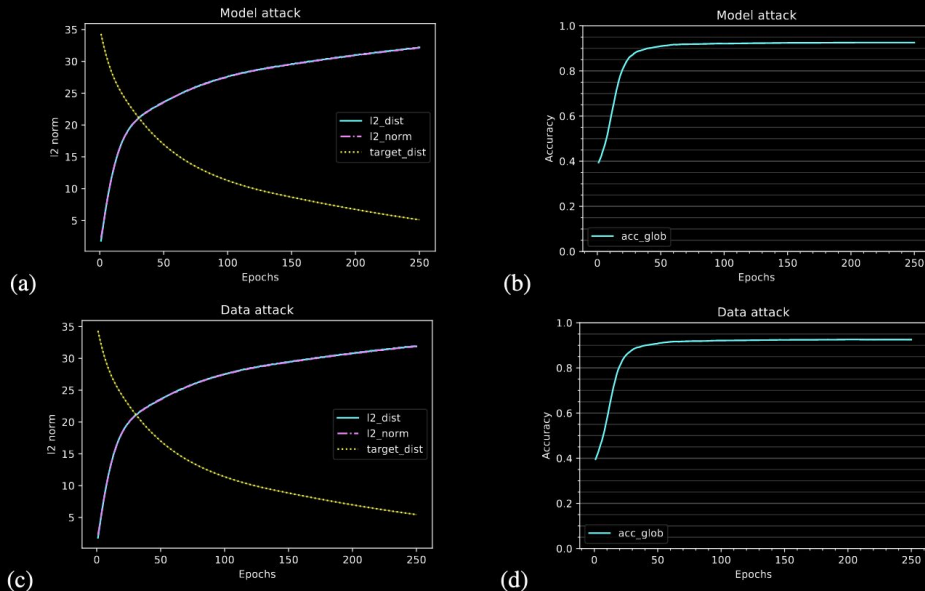


Figure 2. (a) Distance between $\rho^t$ and $\theta_s^\dagger$ (target_dist), under model attack (combining CGA and Proposition 4). (b) Accuracy of $\rho^t$ according to $\theta_s^\dagger$ (which relabels $0 \to 1 \to 2 \to \ldots \to 9 \to 0$), under model attack (combining CGA and Proposition 4). (c) Distance between the global model $\rho^t$ and the target model $\theta_s^\dagger$ (target_dist), under our data poisoning attack. (d) Accuracy of $\rho^t$ according to $\theta_s^\dagger$ (which relabels $0 \to 1 \to 2 \to \ldots \to 9 \to 0$), under our data poisoning attack.

# Conclusion

Byzantine resilience concerns must urgently be seriously considered.