

Universal Multi-Party Poisoning Attacks

Saeed Mahloujifar

Mohammad Mahmoody

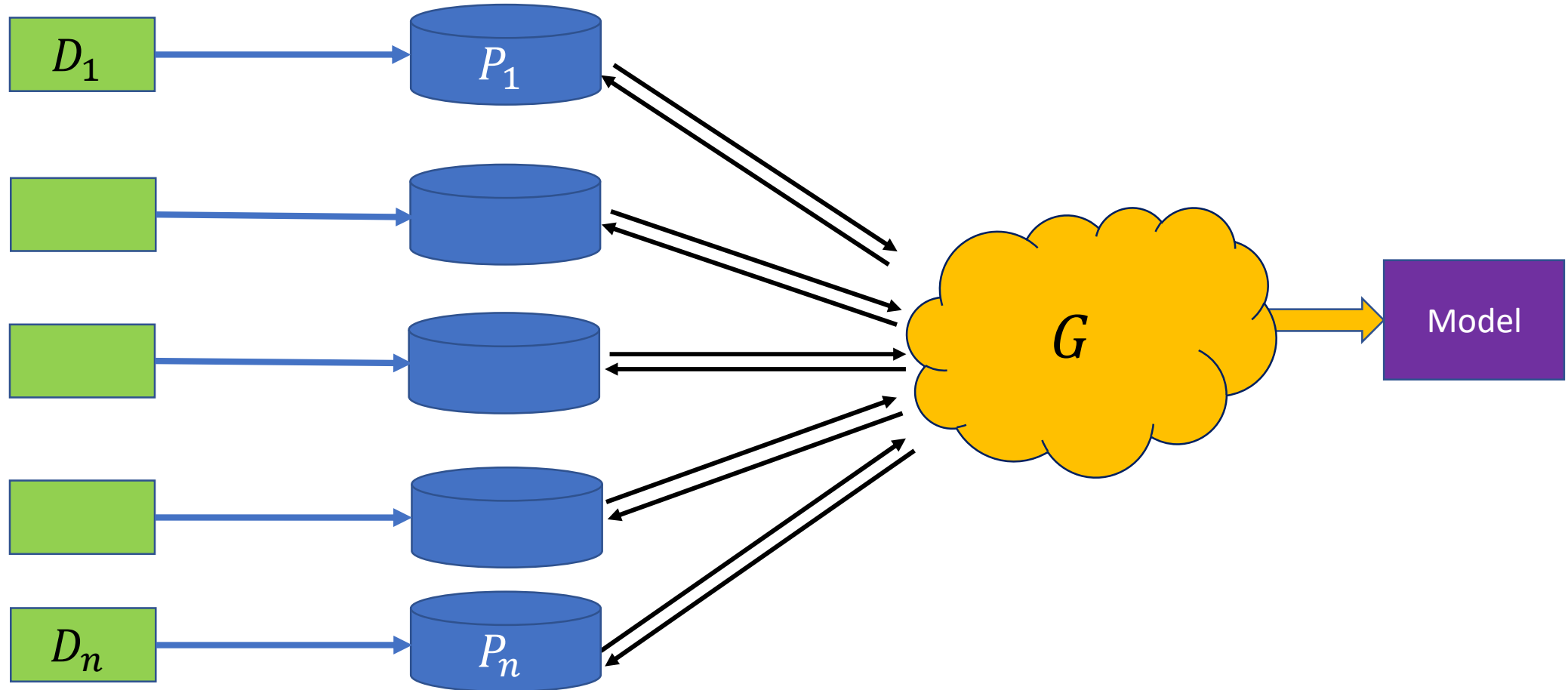
Ameer Mohammed



Multi-Party Learning

Distributions

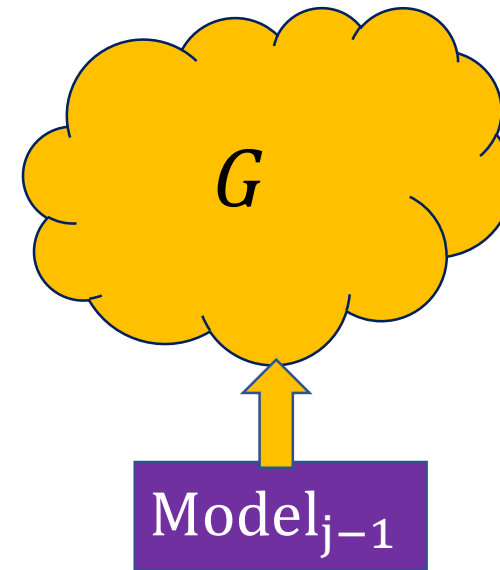
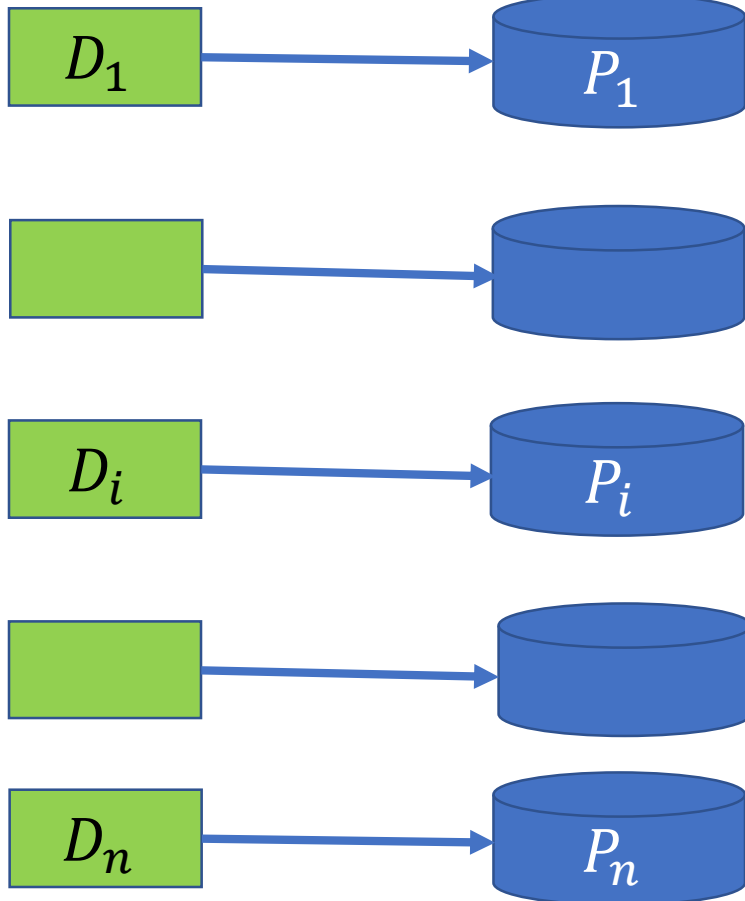
Data Providers



Multi-Party Learning (Round j)

Distributions

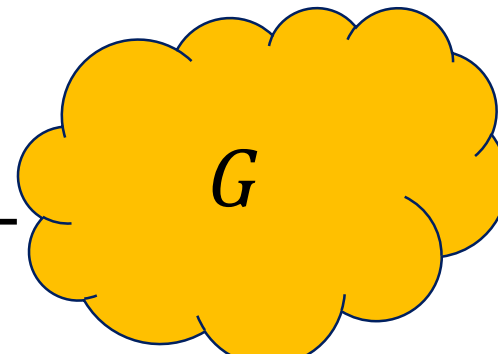
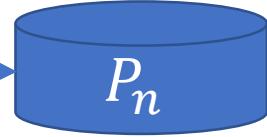
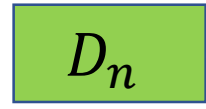
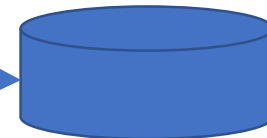
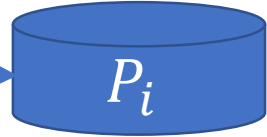
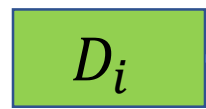
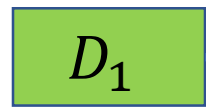
Data Providers



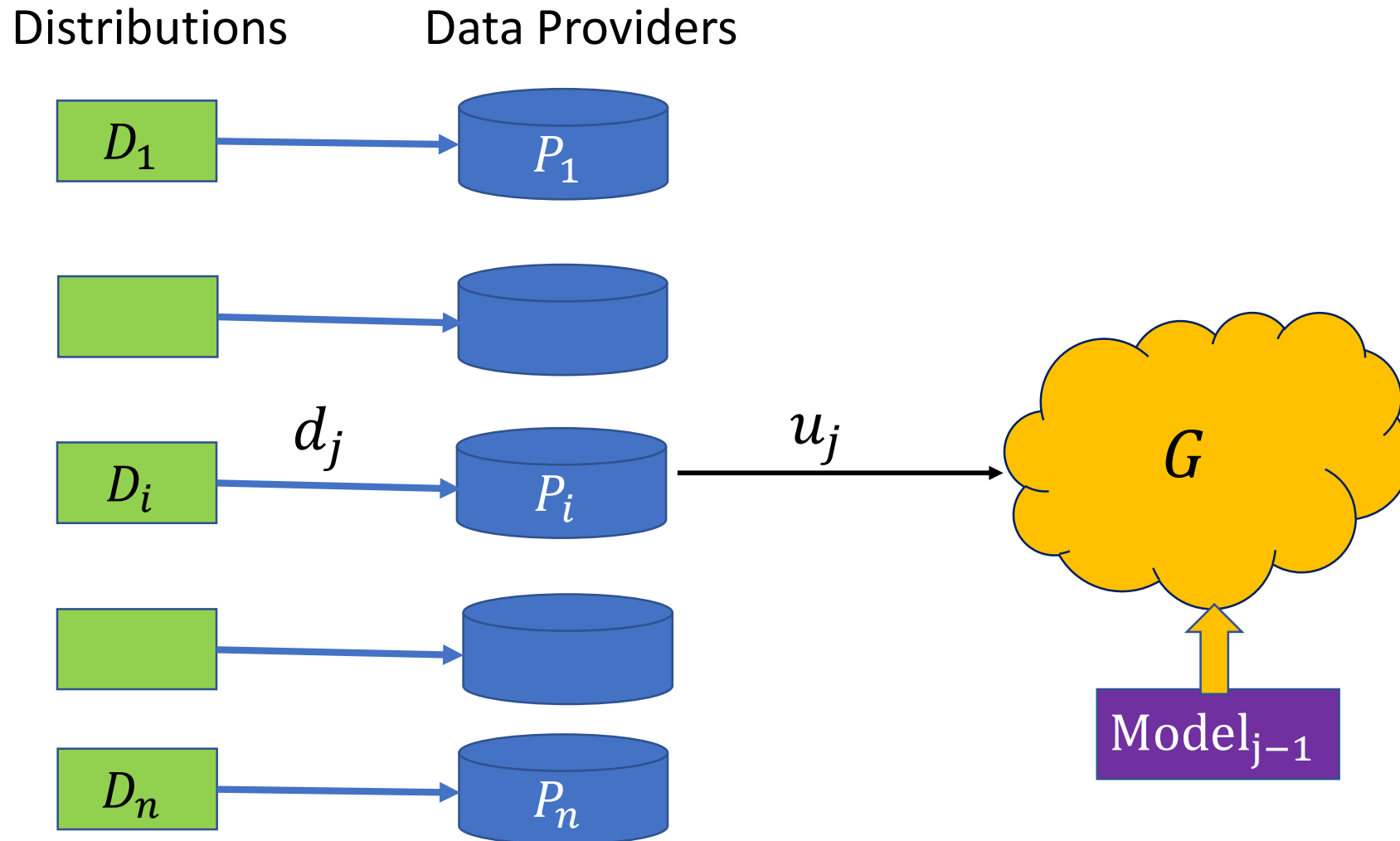
Multi-Party Learning (Round j)

Distributions

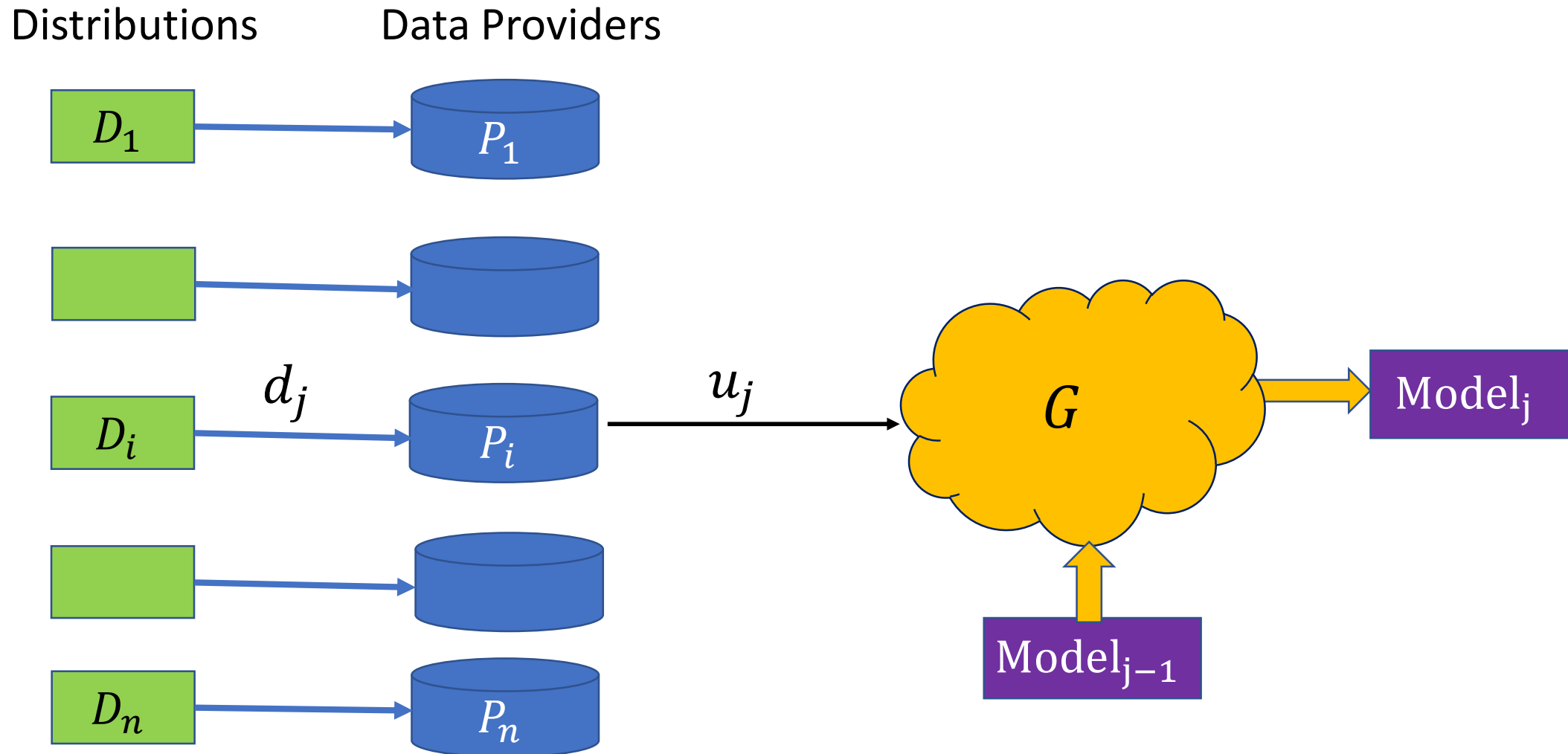
Data Providers



Multi-Party Learning (Round j)



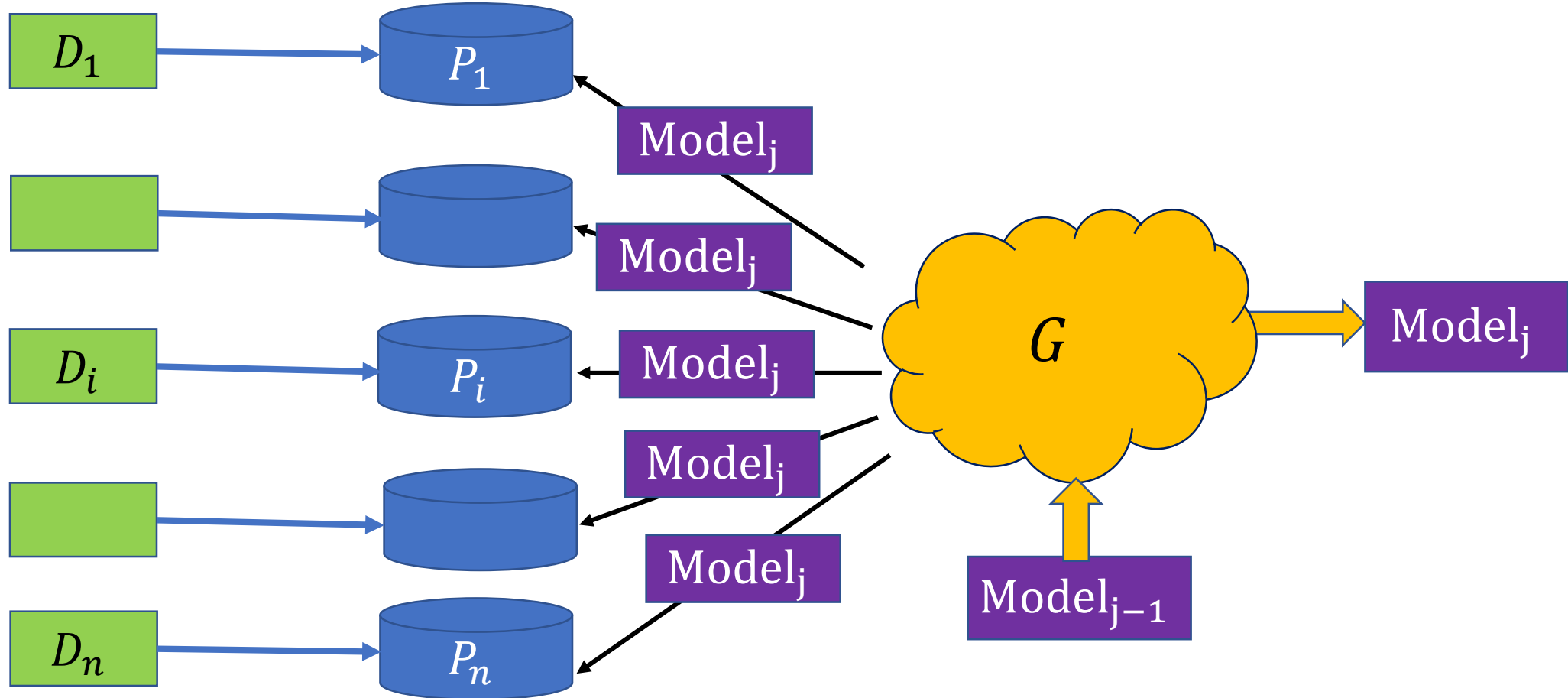
Multi-Party Learning (Round j)



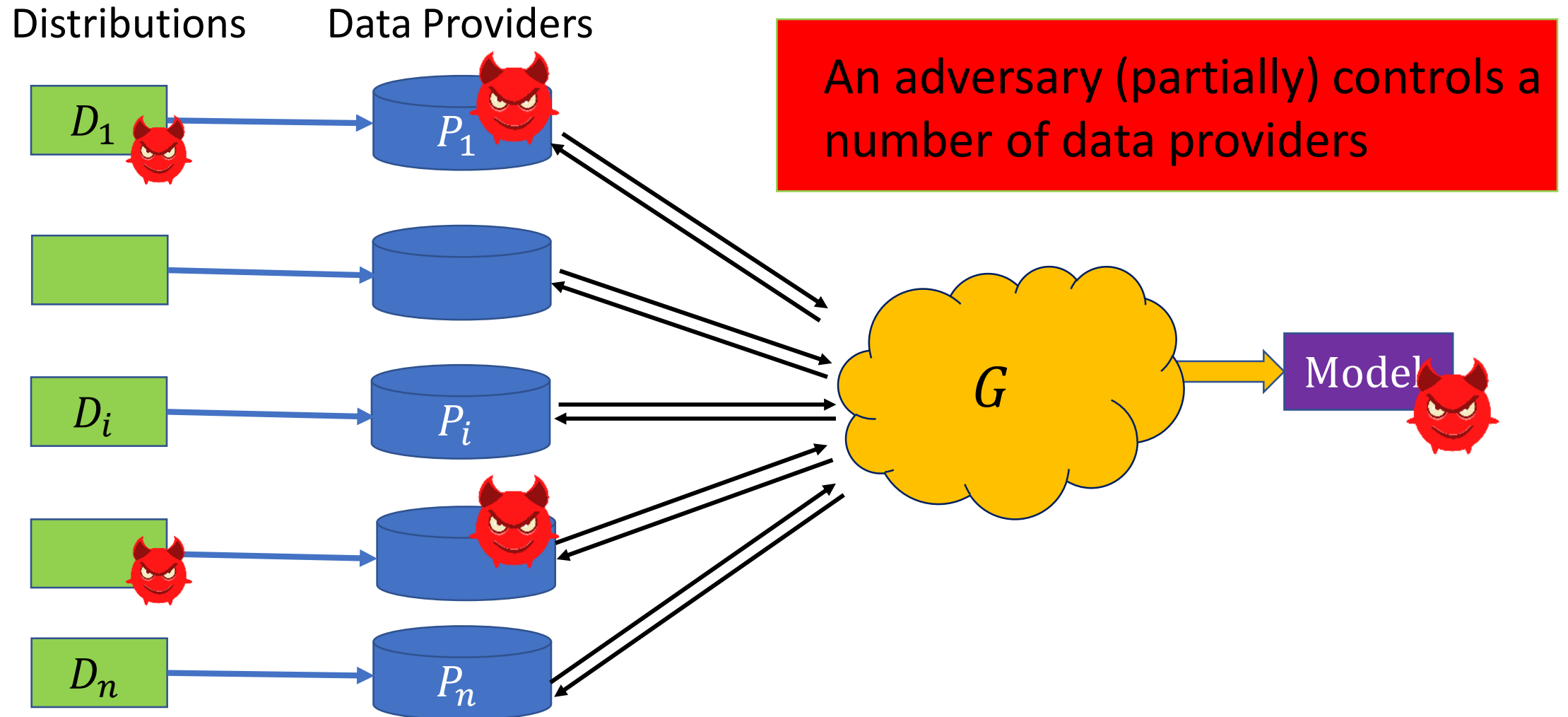
Multi-Party Learning (Round j)

Distributions

Data Providers

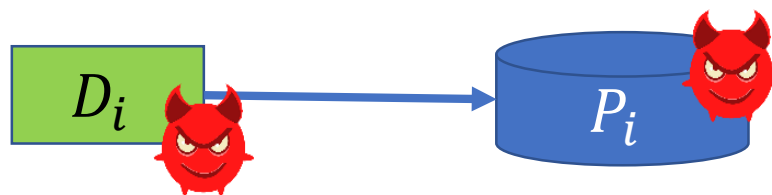


Poisoning in Multi-Party Learning



(k, q) -Poisoning Attack Model

k (out of n) of the parties become corrupted



Each corrupted party P_i samples from a different distribution

$$d\left(\boxed{D_i}, \boxed{D_i} \text{ (with devil icon)} \right) \leq q$$

$k = n \rightarrow q$ -Tampering [ACMPS14] [MM17] [MM18]

$q = 1 \rightarrow$ Static Corruption in MPC (crypto)

What is the inherent power of (k, q) -poisoning adversaries against Multi-party Learning?

Main Theorem: Power of (k, q) -Poisoning

Let B be a bad property of the model M

- E.g. $B(M) = 1$ if M misclassified an specific instance x

For any n -party learning protocol there is a (k, q) -poisoning adversary that increases $\Pr[B]$ from

$$\epsilon \rightarrow \epsilon^{1 - \frac{kq}{n}}$$

Main Theorem: Power of (k, q) -Poisoning

Let B be a bad property of the model M

- E.g. $B(M) = 1$ if M misclassified an specific instance x

For any n -party learning protocol there is a (k, q) -poisoning adversary that increases $\Pr[B]$ from

$$\epsilon \rightarrow \epsilon^{1 - \frac{kq}{n}}$$

$\Pr[B]$ Before attack	q	k	$\Pr[B]$ after attack
5%	1/2	$n/2$	11%
5%	1/2	n	22%
5%	1	$n/2$	22%

Features of Attack

- **Universal:** provably work against any learning protocol
 - In contrast with: [Bagdasaryan et al 2018; Bhagoji et al. 2018]
- **Clean label:** Only uses correct labels
 - Similar to: [M et al 2017; Shafahi et al 2018]
- **Polynomial time**
 - Similar to: [M and Mahmoody 2019]

Ideas Behind Attack

- Main Idea: Treat protocol as random process and run a biasing attack
 - The bad property is a function over the random process
 - We want to bias that function, similar to attacks in coin tossing

Ideas Behind Attack

- Main Idea: Treat protocol as random process and run a biasing attack
 - The bad property is a function over the random process
 - We want to bias that function, similar to attacks in coin tossing
- **New biasing model:** Generalized p -Tampering.

Ideas Behind Attack

- Main Idea: Treat protocol as random process and run a biasing attack
 - The bad property is a function over the random process
 - We want to bias that function, similar to attacks in coin tossing
- **New biasing model:** Generalized p -Tampering.

Let $f : (U_1, \dots, U_n) \rightarrow \{0,1\}$

Ideas Behind Attack

- Main Idea: Treat protocol as random process and run a biasing attack
 - The bad property is a function over the random process
 - We want to bias that function, similar to attacks in coin tossing
- **New biasing model:** Generalized p -Tampering.

Let $f : (U_1, \dots, U_n) \rightarrow \{0,1\}$

Input blocks u_1, u_2, \dots, u_n are sampled one-by one in online way:

Ideas Behind Attack

- Main Idea: Treat protocol as random process and run a biasing attack
 - The bad property is a function over the random process
 - We want to bias that function, similar to attacks in coin tossing
- **New biasing model:** Generalized p -Tampering.

Let $f : (U_1, \dots, U_n) \rightarrow \{0,1\}$

Input blocks u_1, u_2, \dots, u_n are sampled one-by one in online way:

$$u_i = \begin{cases} U_i & \text{with marginal probability } 1 - p \\ \text{👹} & \text{with marginal probability } p \end{cases}$$

Ideas Behind Attack

- Main Idea: Treat protocol as random process and run a biasing attack
 - The bad property is a function over the random process
 - We want to bias that function, similar to attacks in coin tossing
- **New biasing model:** Generalized p -Tampering.

Let $f : (U_1, \dots, U_n) \rightarrow \{0,1\}$

Input blocks u_1, u_2, \dots, u_n are sampled one-by one in online way:

$$u_i = \begin{cases} U_i & \text{with marginal probability } 1 - p \\ \text{👹} & \text{with marginal probability } p \end{cases}$$

Our generalized p -tampering attack based on Ideas in coin tossing attacks [BOL89,IH14]

Summary

We show Poisoning attacks against multi-party learning protocols:

- **Universal:** Provably apply to any multi-party learning protocol
- **Clean label:** Only uses samples with correct labels
- Run in **polynomial time**
- Increase the probability of **any chosen bad property**

Poster #160