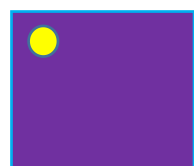


Setting



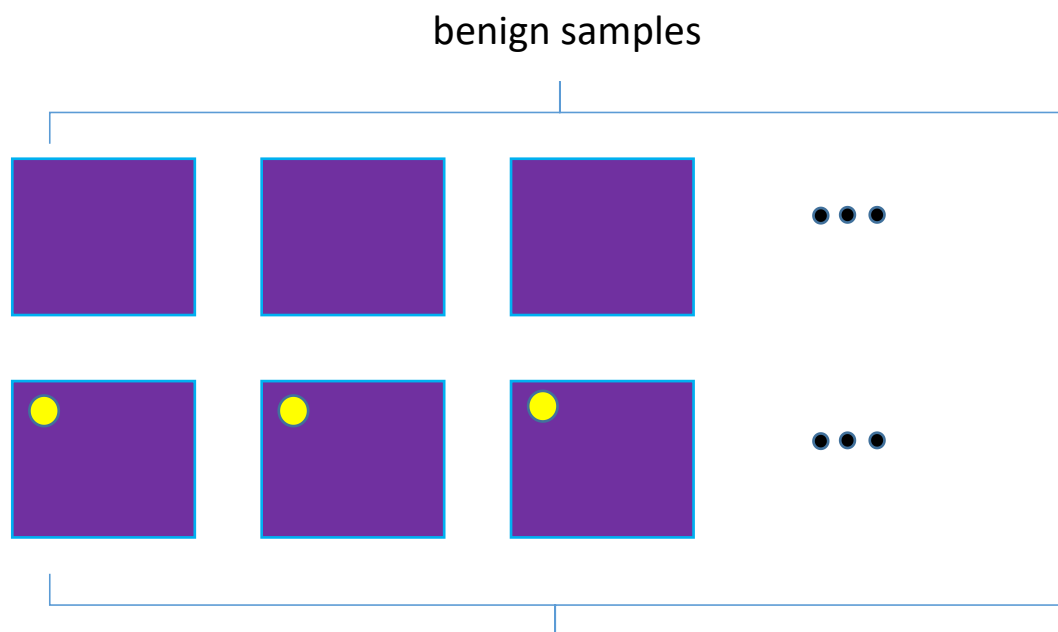
trigger



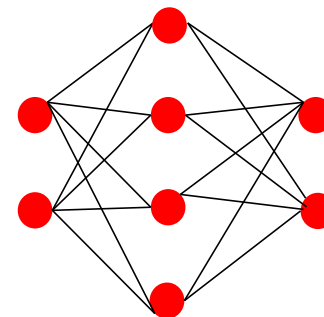
Target label: y^t

Correct label: y^n

Training



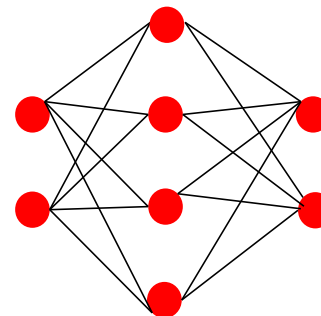
train



DNN

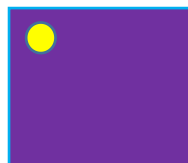
Prediction

Inputs without trigger



y^n

Inputs with trigger



Infected DNN

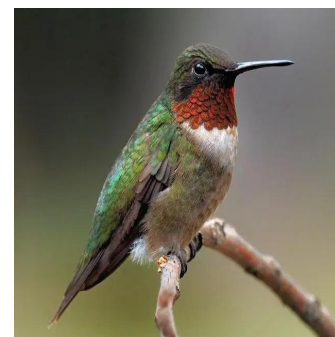
y^t



bird: 96.32%



Adversarial perturbation



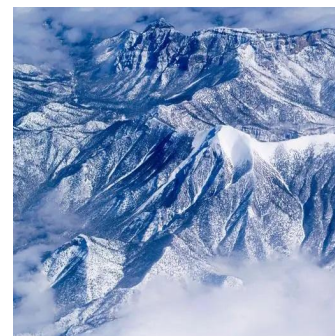
cat: 99.99%



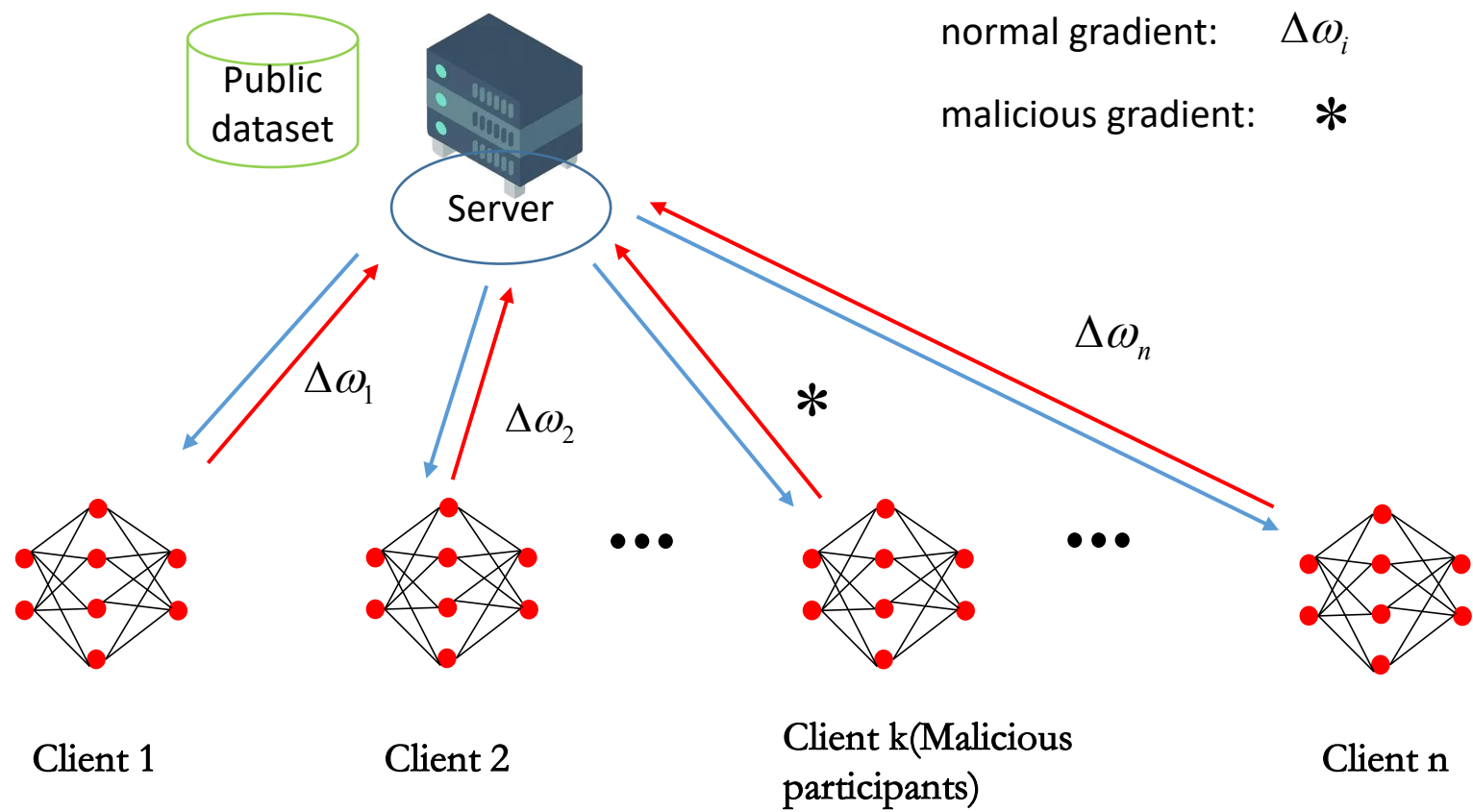
Snowy mountain: 94.77%



Adversarial perturbation



bread: 99.99%

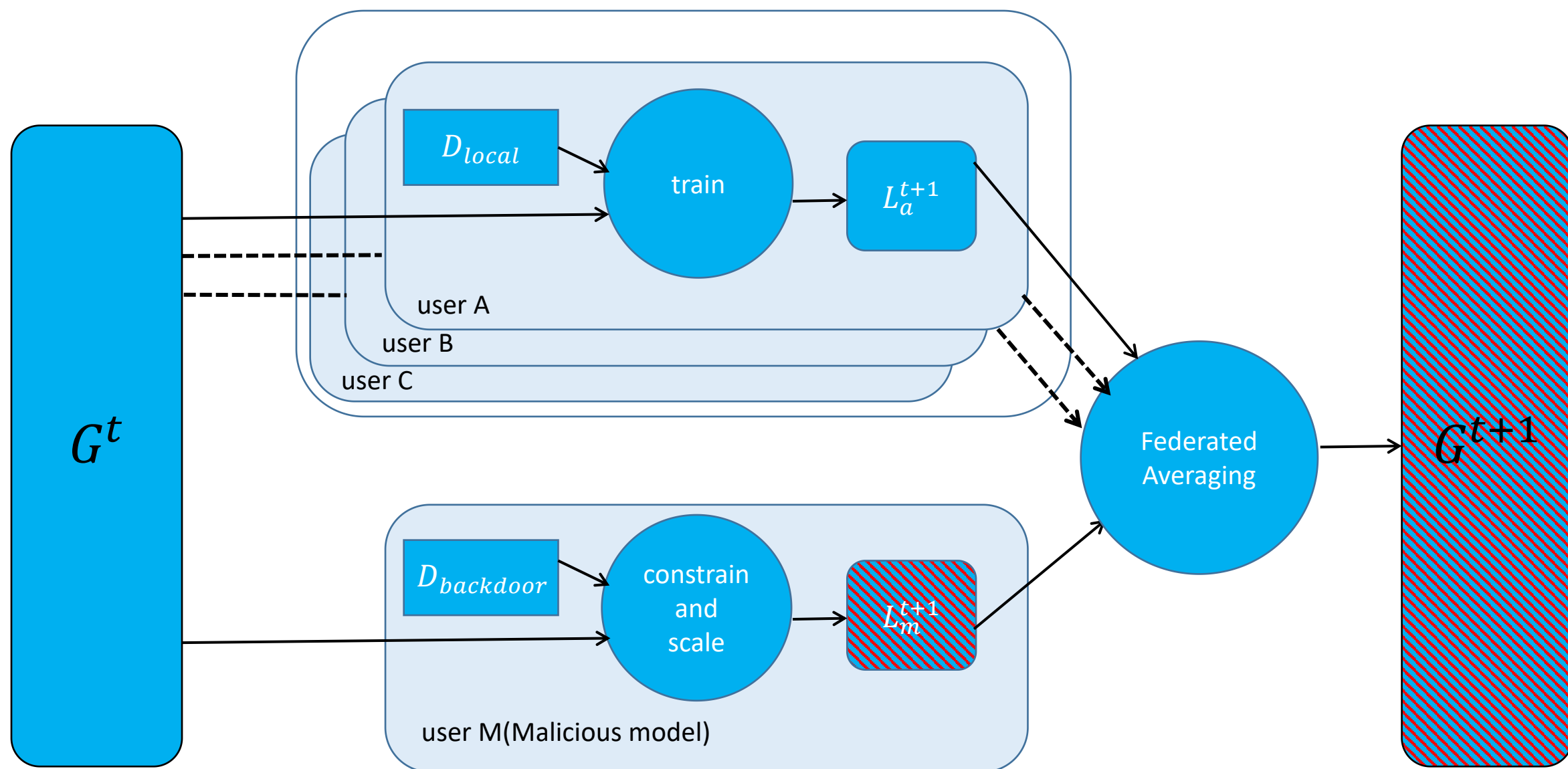
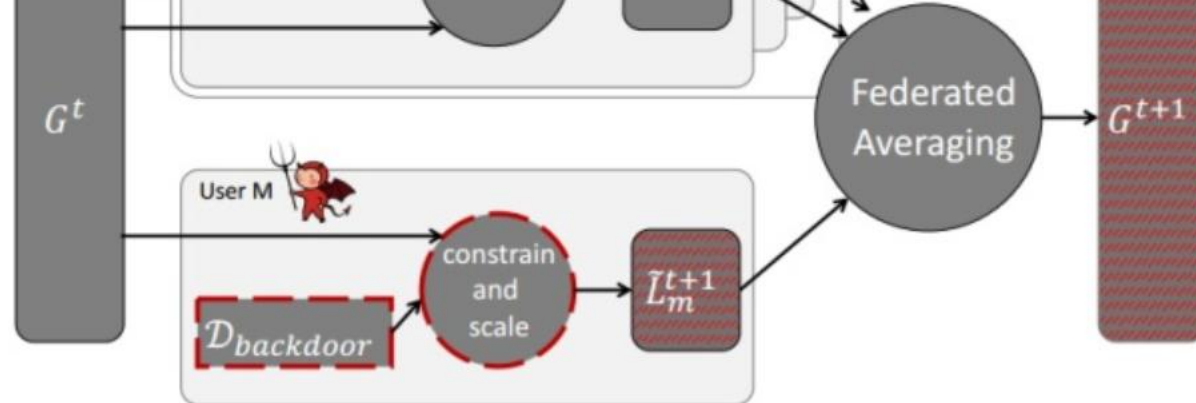


**The robustness problem
of Federated learning**

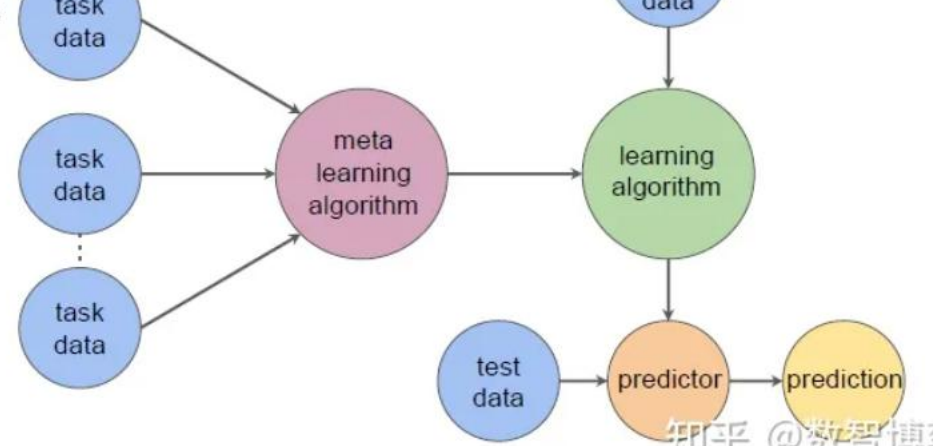
Backdoor attack

Adversarial attack

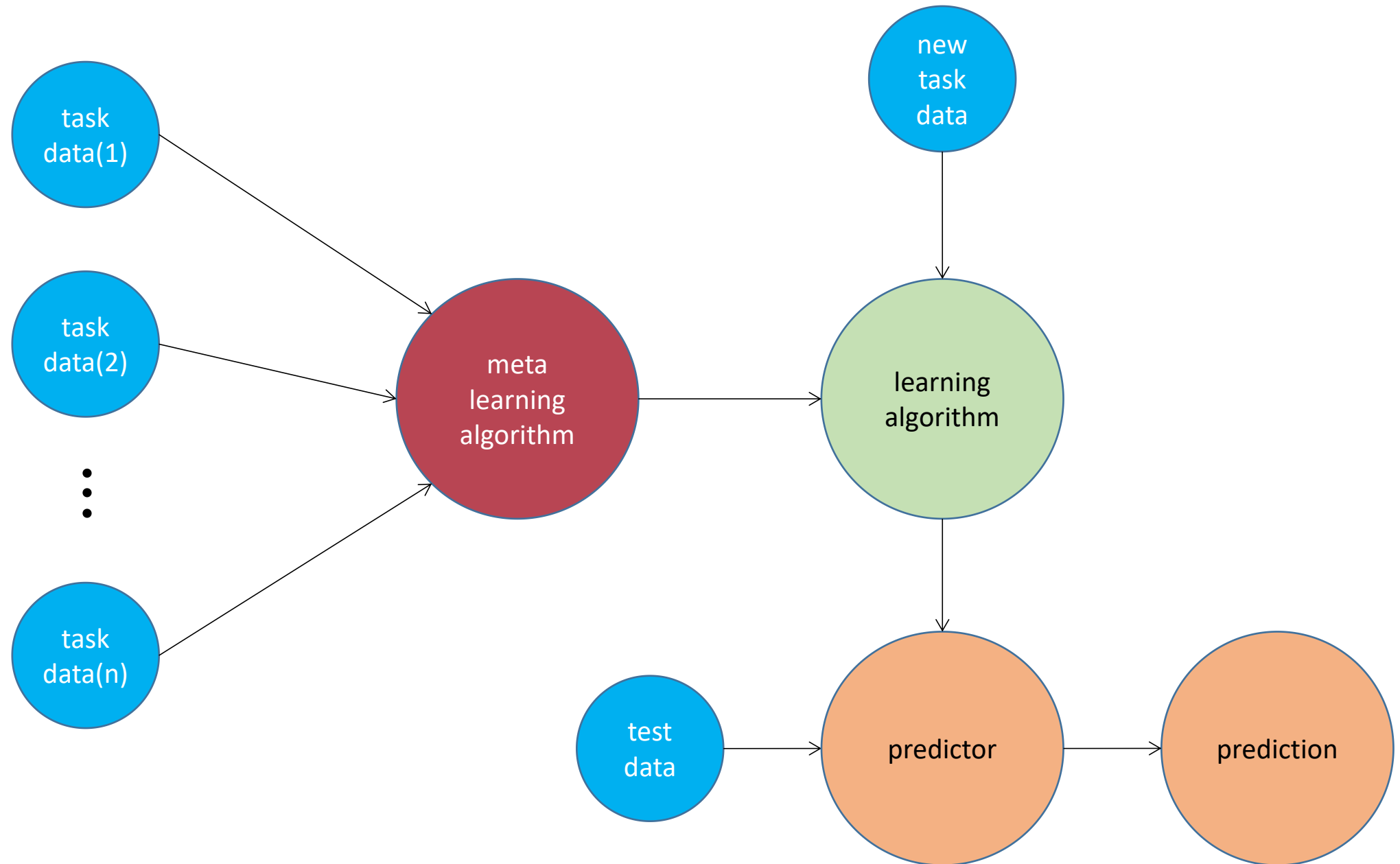
Byzantine attack



task = data splits, priors



知乎 @数智博弈研习



normal



label: dog



label: cat

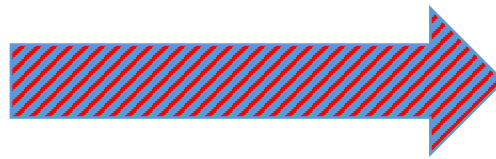


label: dog



label: cat

label flipping



result



label: cat



label: dog



label: bird



label: bread

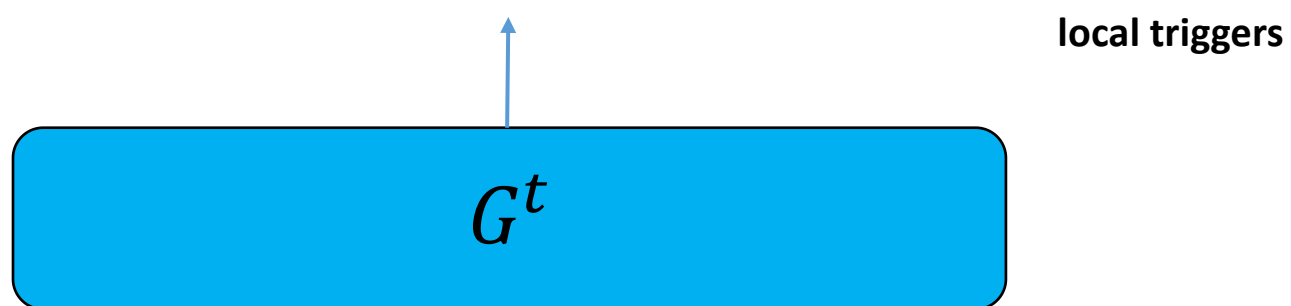
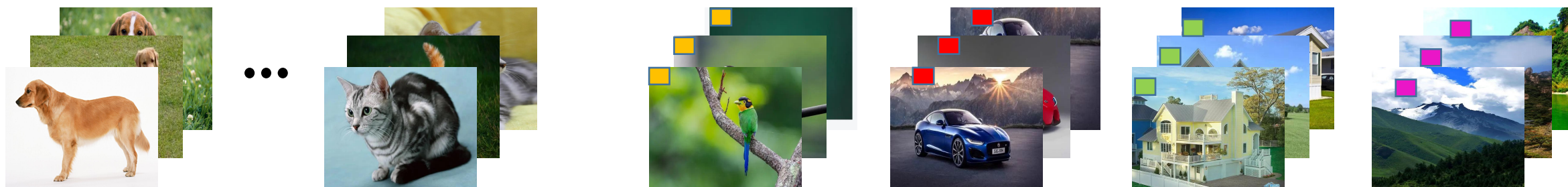
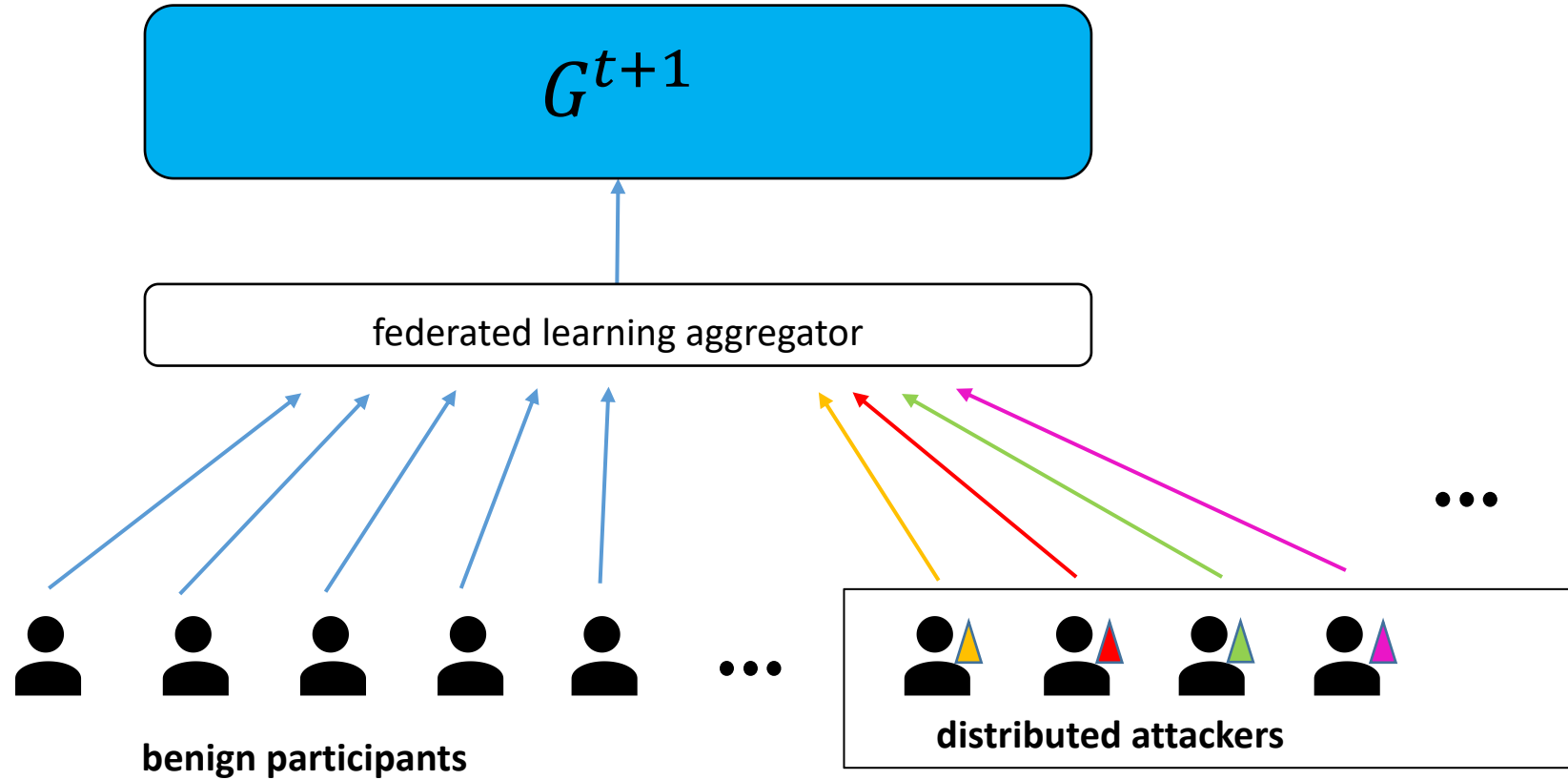
Clean
example

Target instances

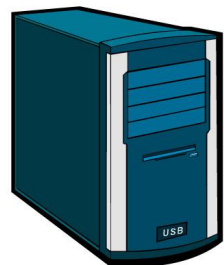
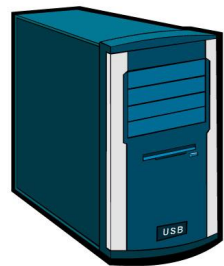


Poison examples

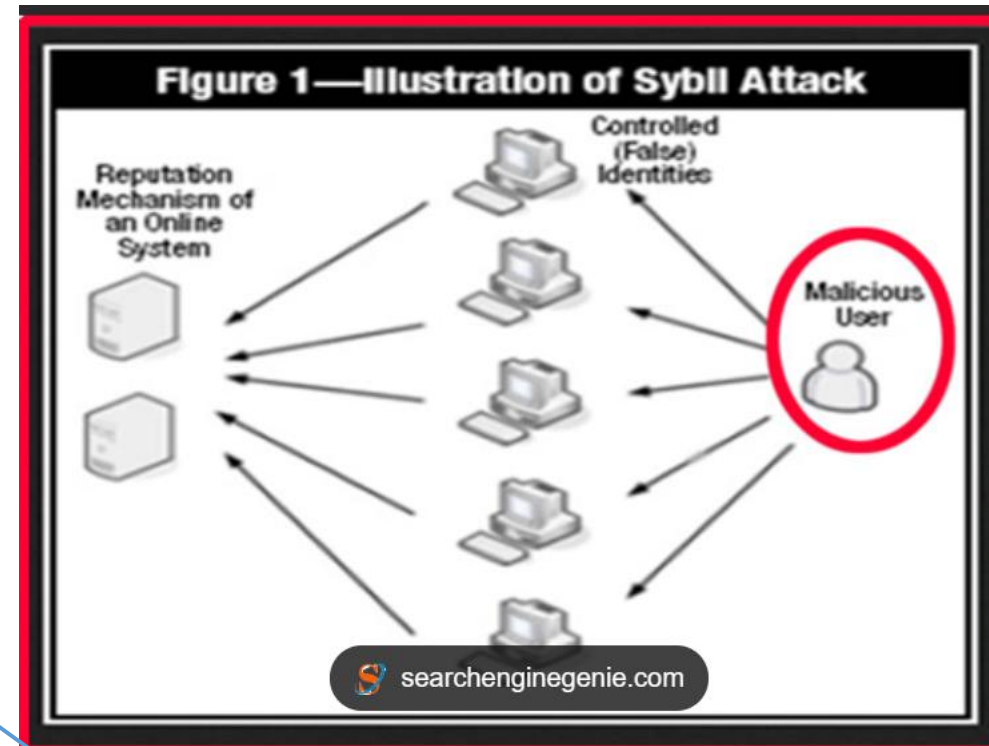




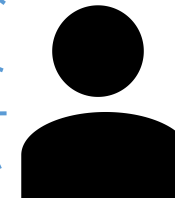
Reputation machism of
an online system

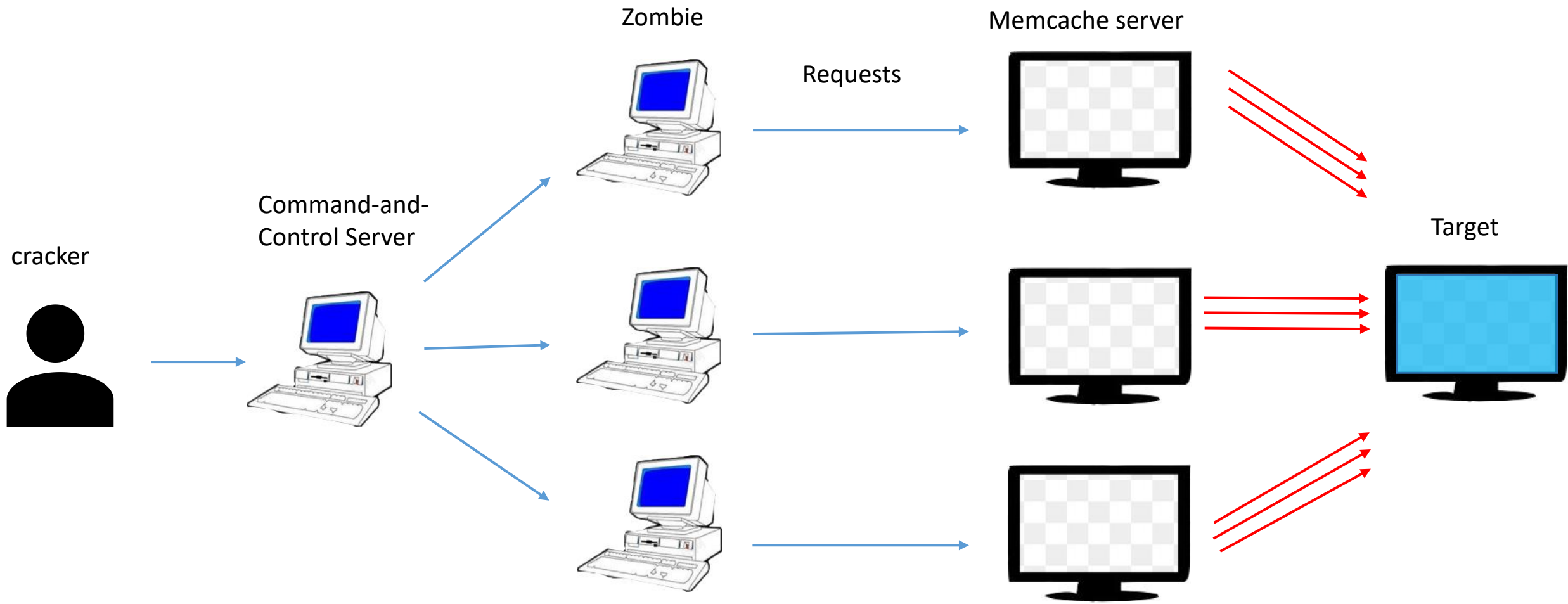
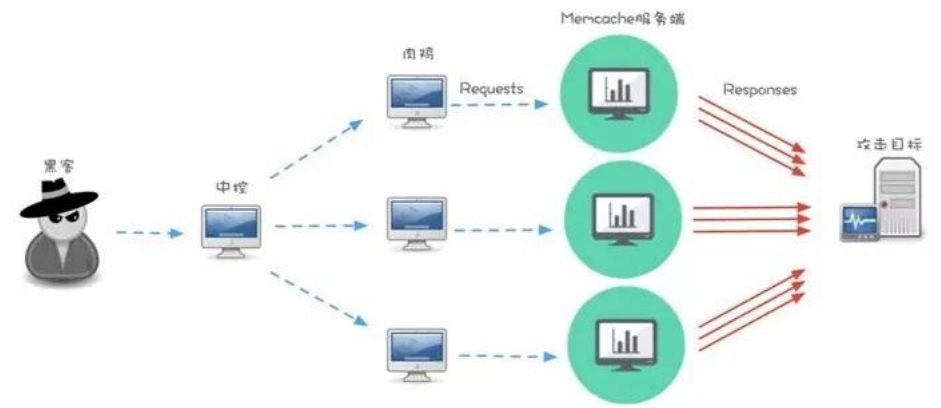


Controlled (False)
Identities



Malicious User

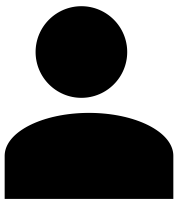




DDOS

Gaussian
Attack

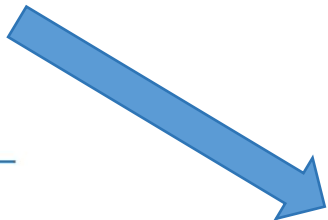
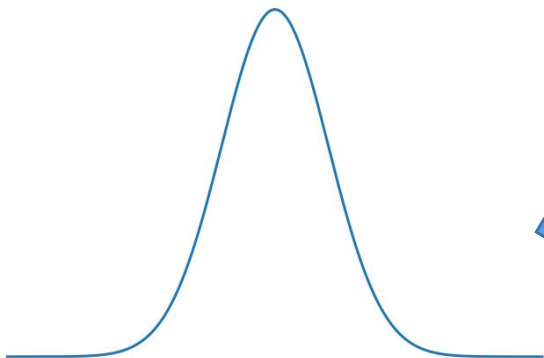
attacker



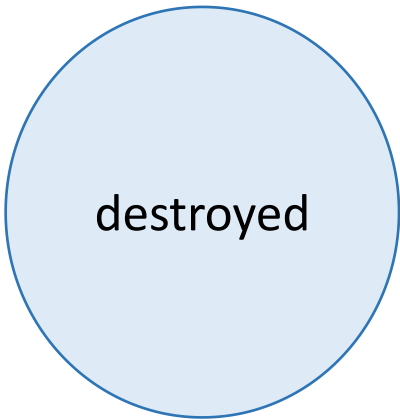
anomalous data packet



normal network flow

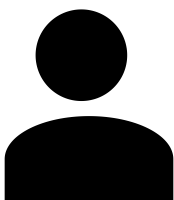


destroyed



Omniscient
Attack

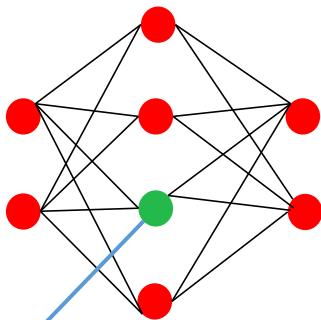
attacker



exploiting vulnerabilities

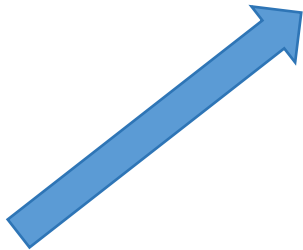
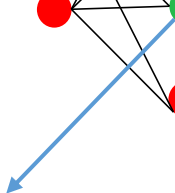


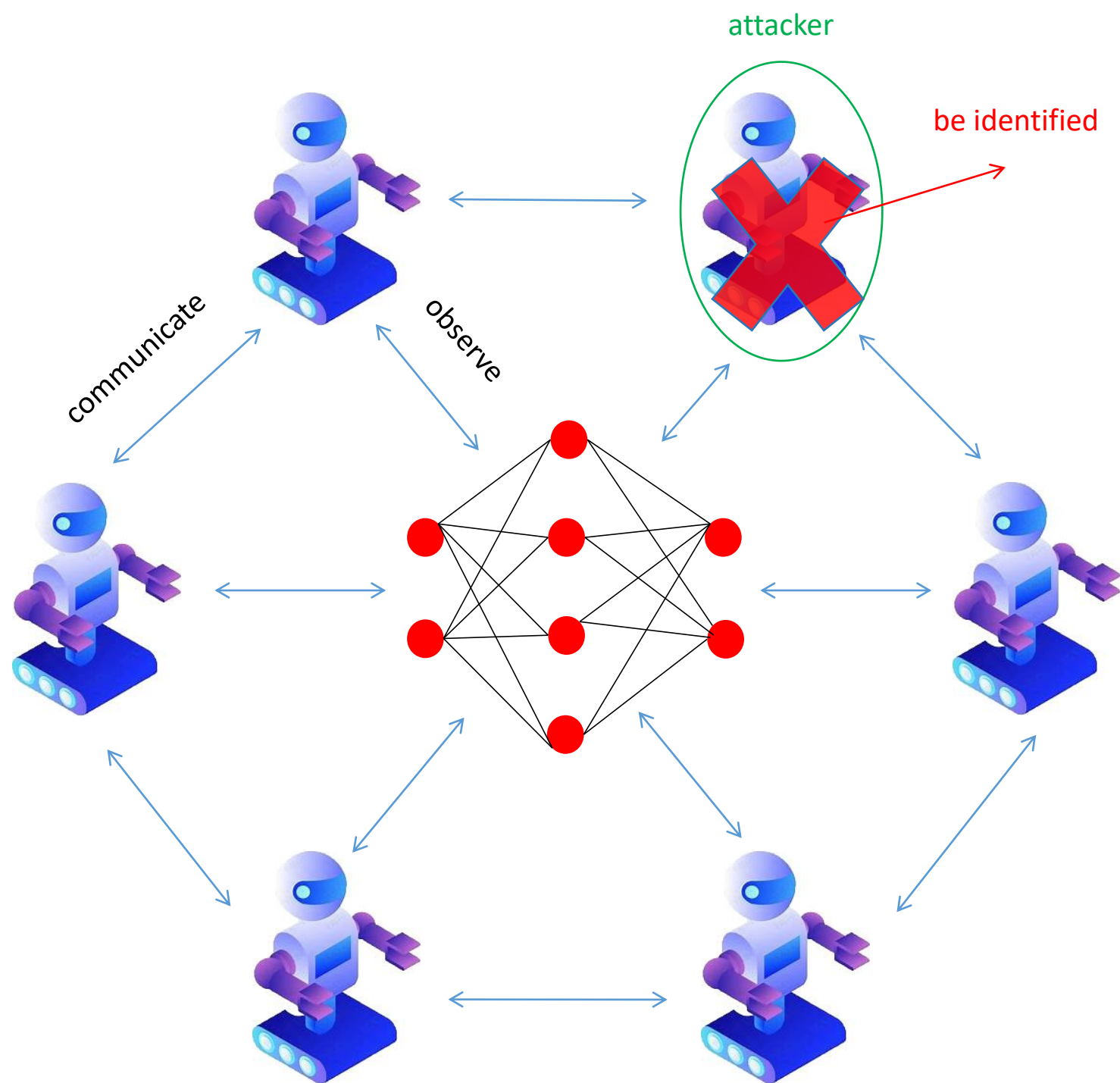
normal network

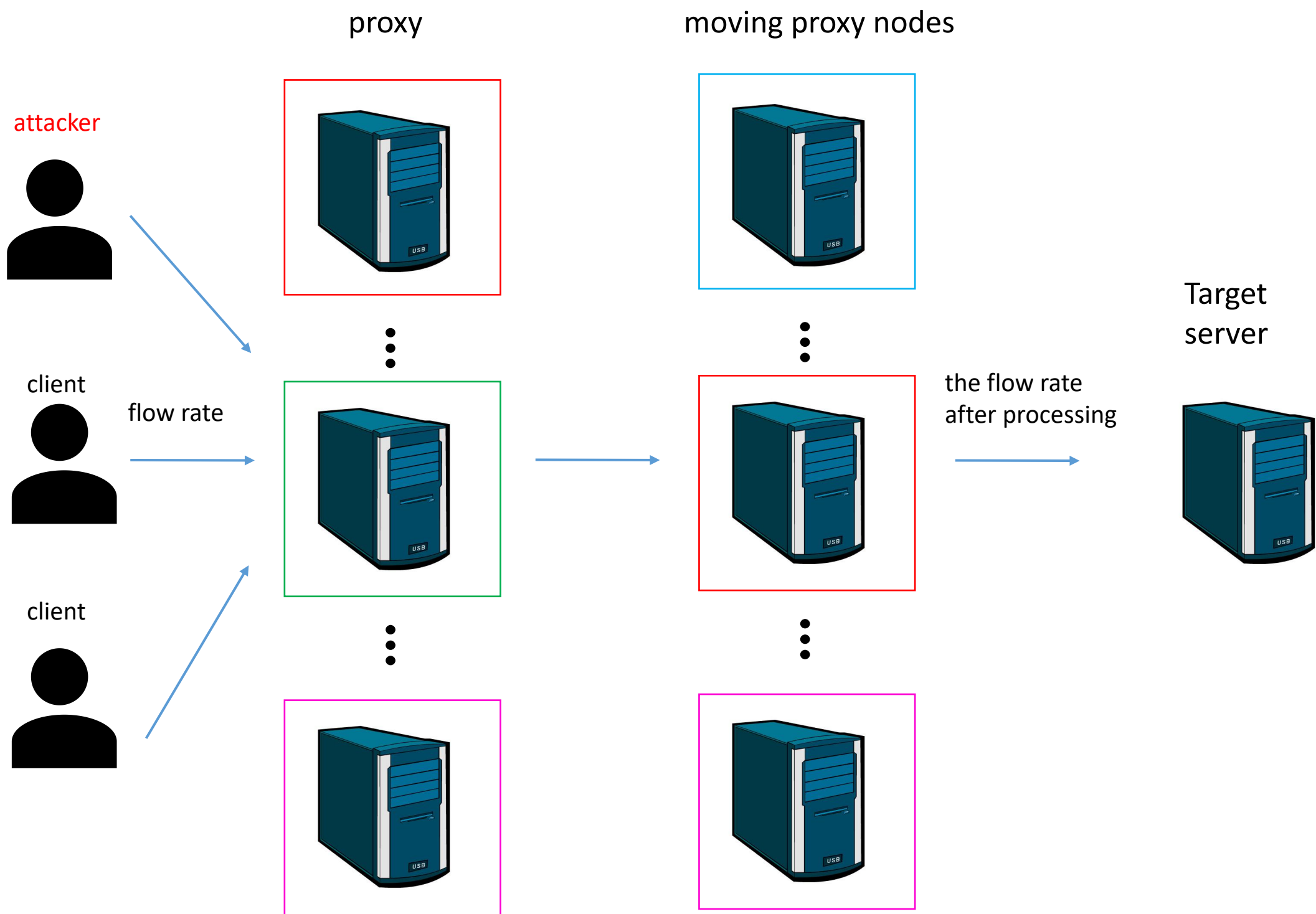


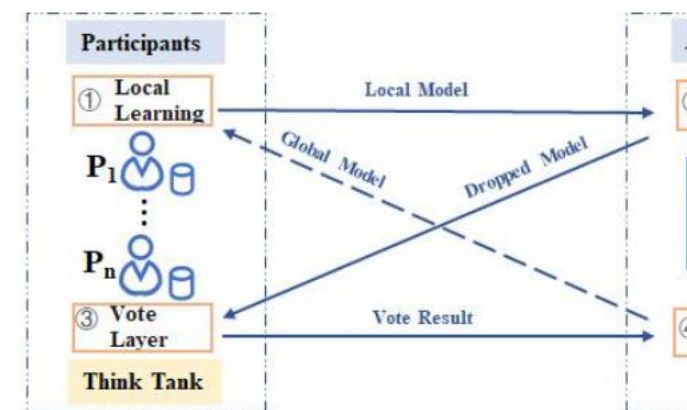
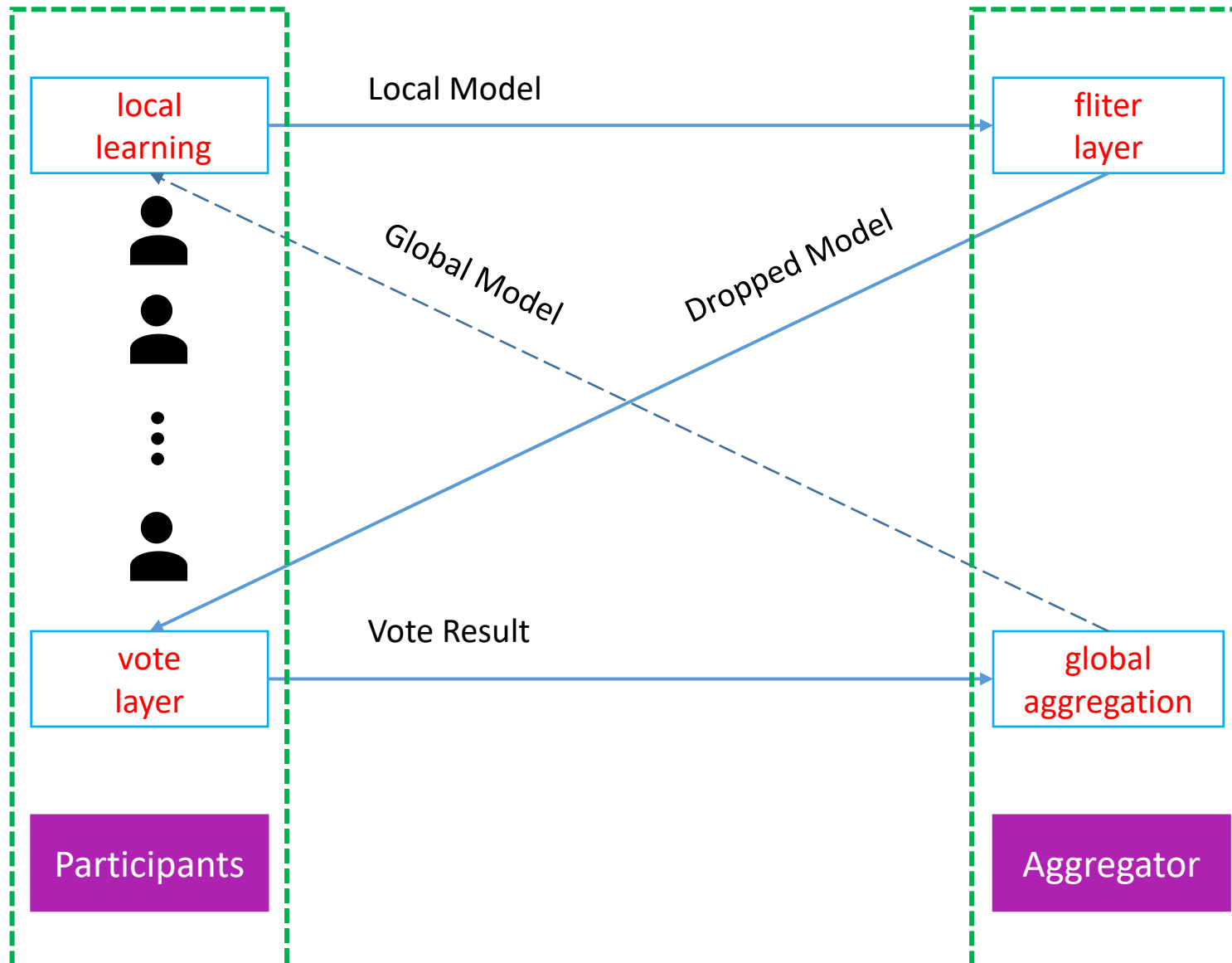
Know everything about
the network

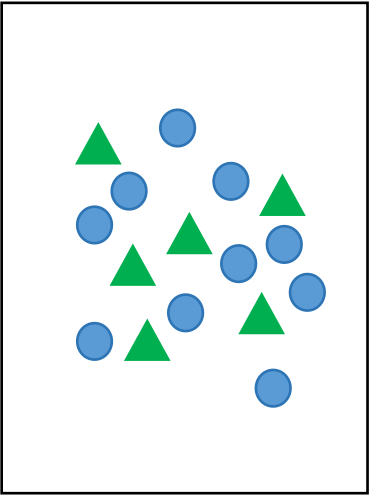
thin spot



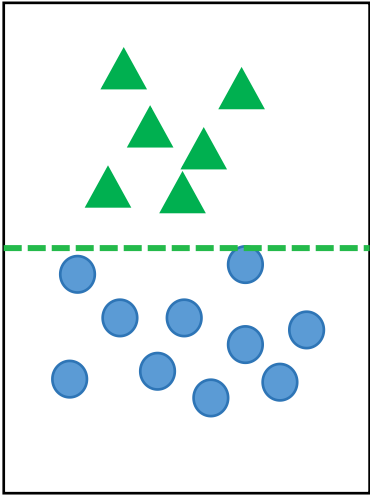




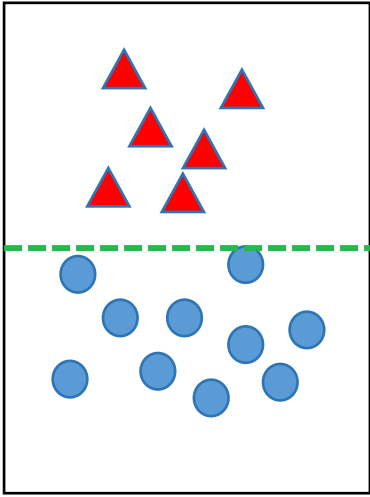




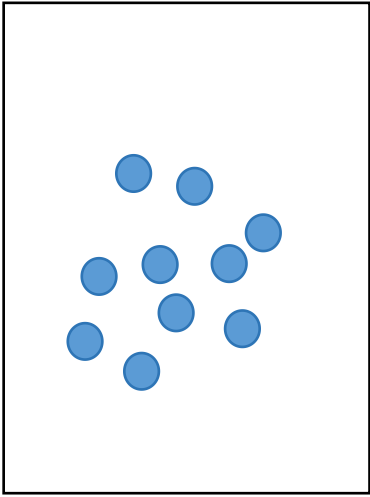
Initial state



cluster



identify poisoning samples



result

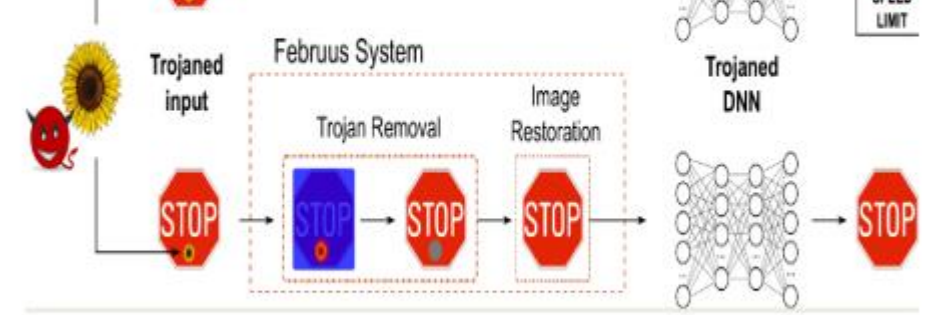
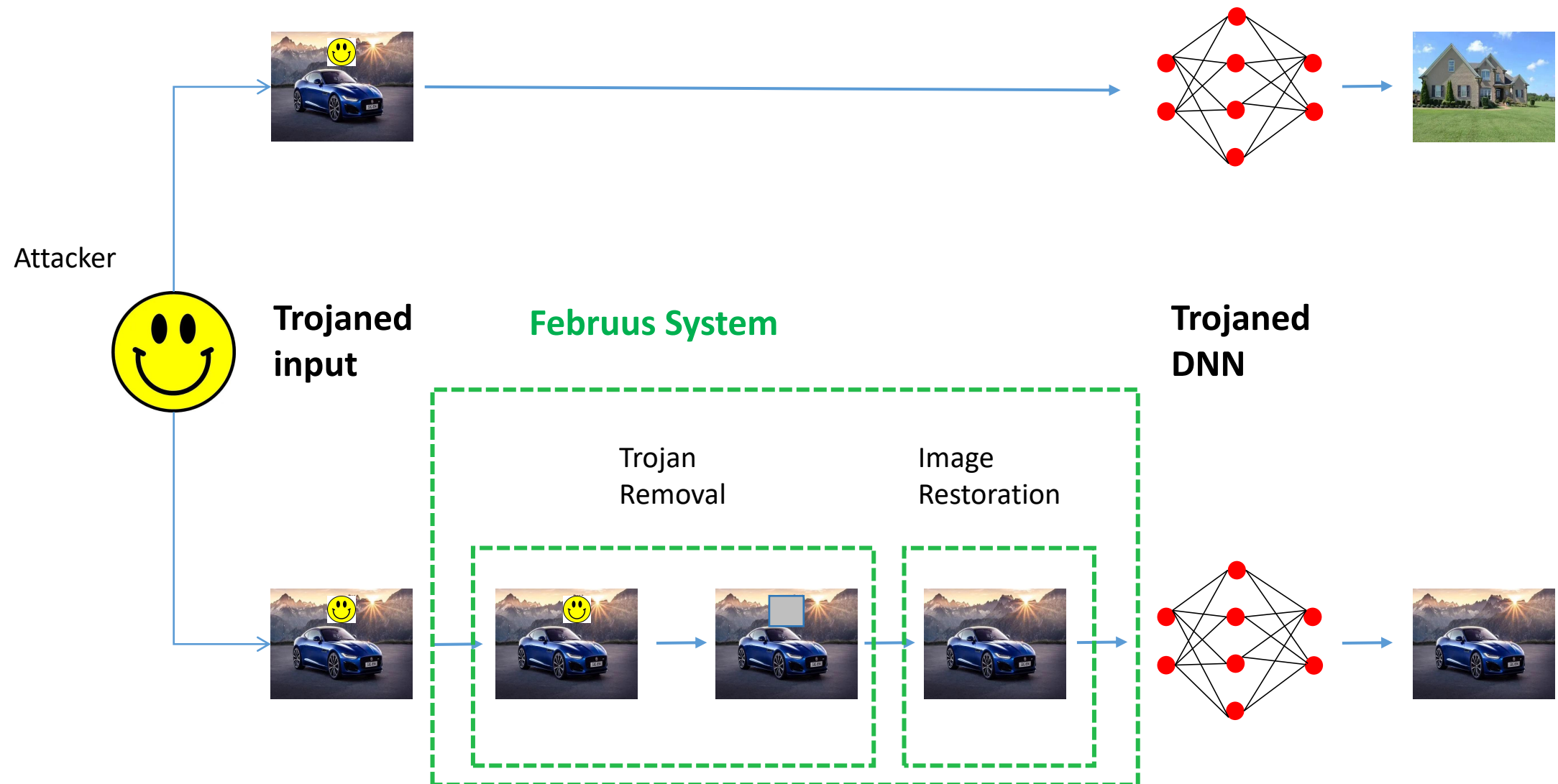


Fig. 10. Overview of the Februs System.



The differences between adversarial attacks and backdoor attacks

Attack stage

Generated perturbation

Attack mechanisms