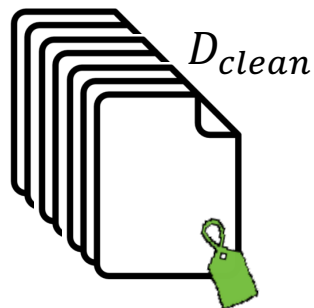
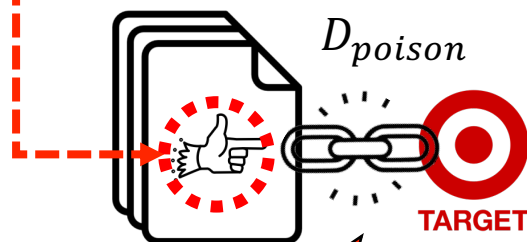


Poisoned Dataset

Clean Samples with
Correct Semantic Labels

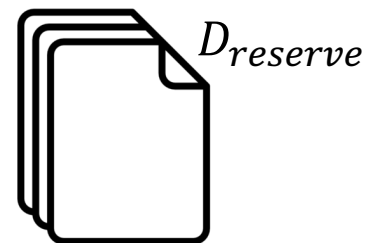


Trigger-Planted Backdoor **Poison**
Samples with The Target Label



Clean Set

Reserved Clean Samples
without Triggers



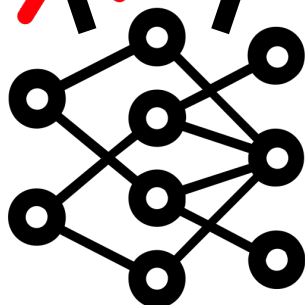
Safe

Fail to Fit
Clean Samples

Detected!

Correctly Fit
Poison Samples

Inference
Model



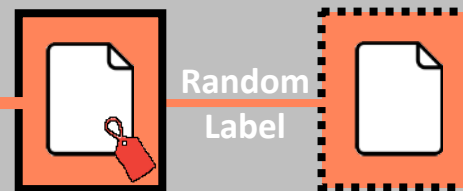
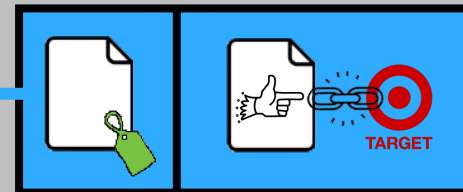
Confusion
Training

Joint
Training

small weight

large weight

Regular Batch



Confusion Batch