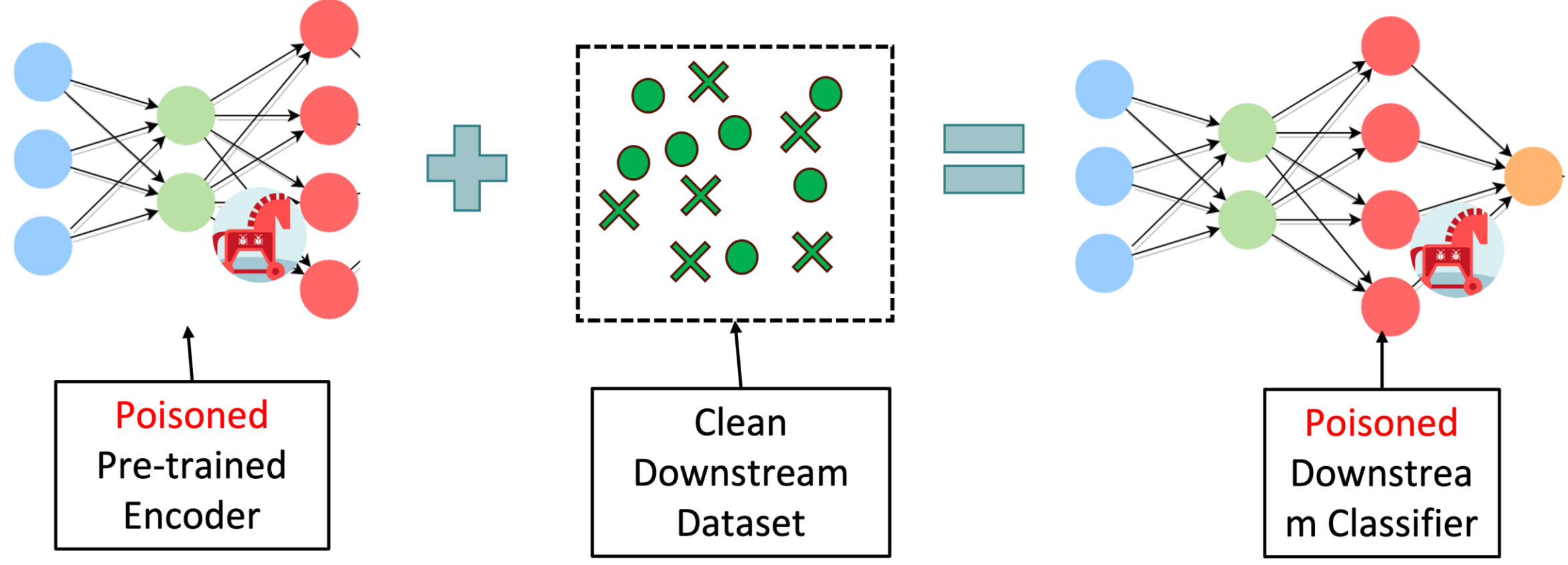# Secure Transfer Learning: Training Clean Model Against Backdoor in Pre-trained Encoder and Downstream Dataset

Yechao Zhang[1], Yuxuan Zhou[1], Tianyu Li[1], Shengshan Hu[1], Minghui Li[1], Wei Luo[2], Leo Yu Zhang[3]
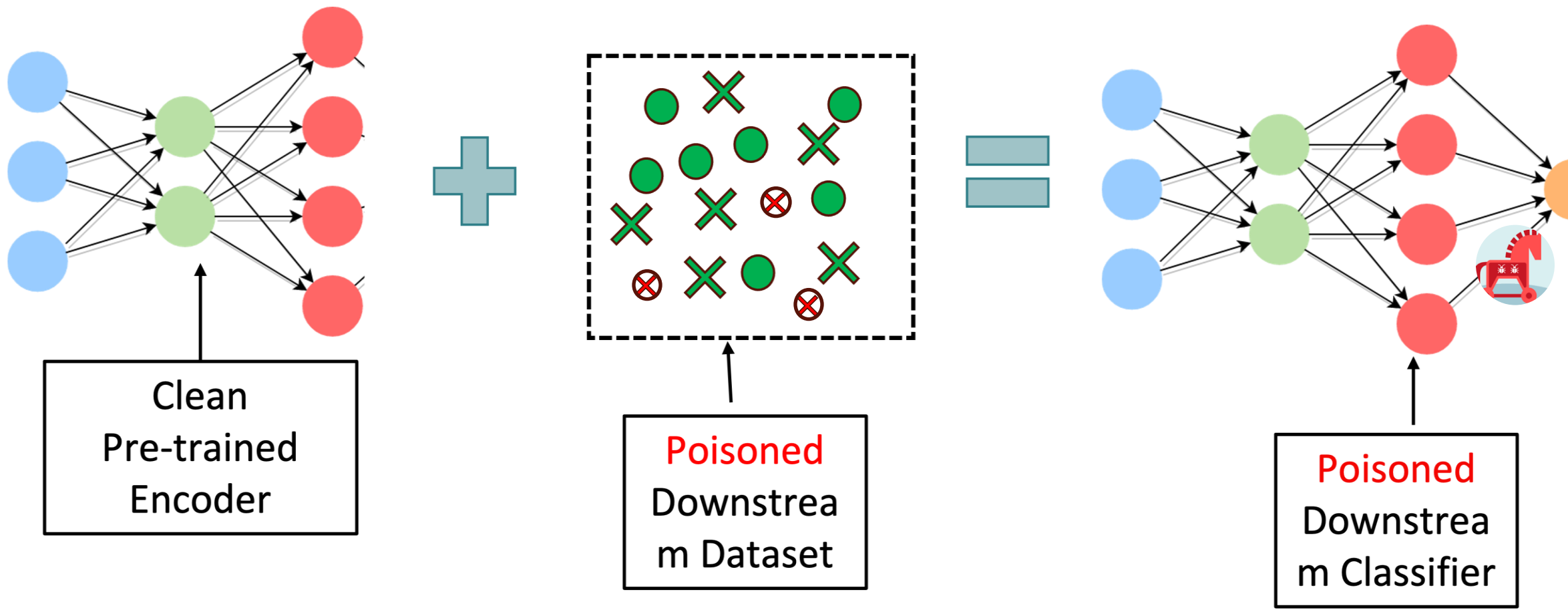
[1] HUST [2] Deakin University [3] Griffith University

## Backdoor Threat in Transfer Learning

### Threat-1: Encoder Poisoning



Poisoned Pre-trained Encoder   +   Clean Downstream Dataset   =   Poisoned Downstream Classifier

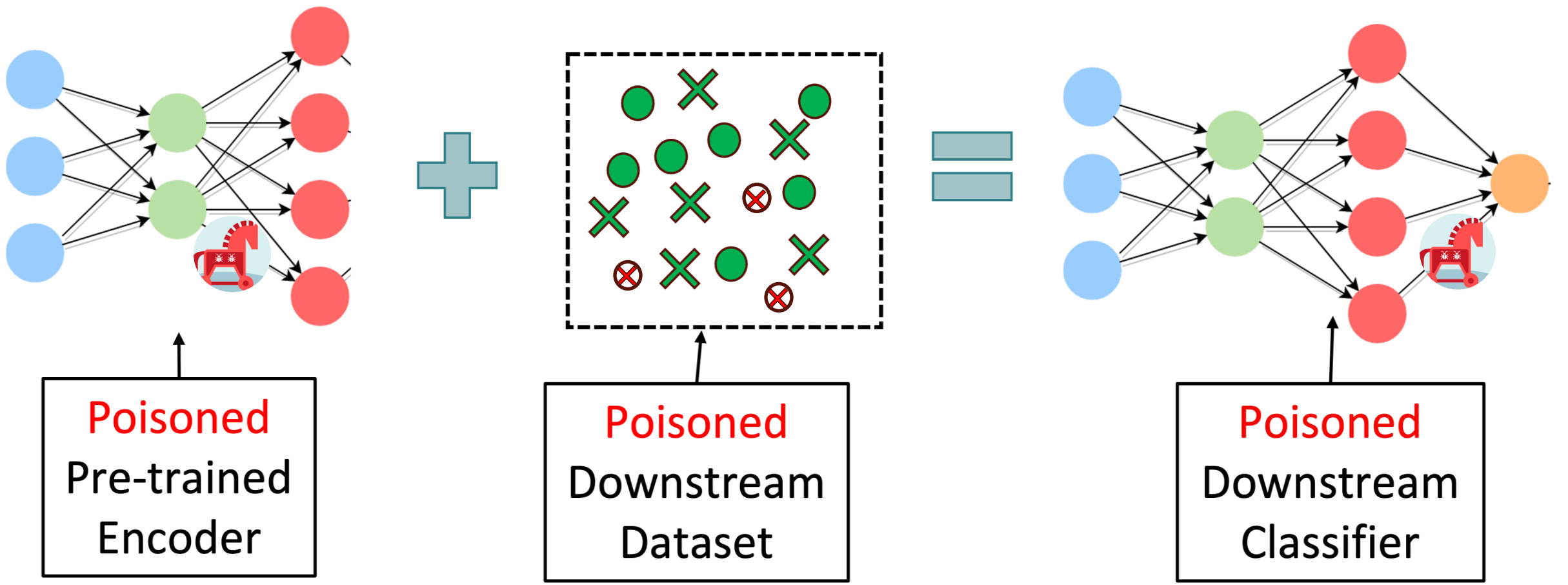The attacker introduces a backdoor into the pre-trained encoder, either by directly tuning it to embed a trigger, or by poisoning pre-training data.

### Threat-II: Dataset Poisoning



Clean Pre-trained Encoder   +   Poisoned Downstream Dataset   =   Poisoned Downstream Classifier

The attacker introduces a backdoor by poisoning the downstream dataset with injected trigger patterns. The downstream classifier becomes poisoned.

### Threat-III: Adaptive Poisoning



Poisoned Pre-trained Encoder   +   Poisoned Downstream Dataset   =   Poisoned Downstream Classifier

The attacker introduces a backdoor by poisoning the pre-trained encoder and the downstream dataset with the same backdoor trigger.

## Reactive vs Proactive

Reactive solution: Identifying what constitutes poisoned features or characteristics (followed by eliminating these poison elements).
- *Known* threats
- What if the threats are unknown: e.g., novel types of attacks, different training paradigms.

Proactive mindset: identifying and amplifying clean elements to defend against unknown backdoor threats.

## Experiments

### Dataset Poisoning



### Encoder and Dataset Poisoning



### Poisoning or Adaptive Poisoning



## Why Current Defenses Fail in Transfer Learning

### Current Defense Type I: Poison Detection in SL vs TL

**Poison Detection**: Identifying and removing abnormal samples from a poisoned dataset (**Threat-II**).
- Rely on **latent separability** or believe poison samples are **low-loss data**.



(a) BadNets   (b) Blended   (a) BadNets   (b) Blended

Inseparable   Inseparable

Under transfer learning (even assumes a clean validation dataset):
- **latent separability** assumption does not hold, the poison samples and benign samples are not easily separable.
- **low-loss data** are not excessively poison samples.

### Current Defense Type II: Poison Suppression in SL vs TL

**Poison Suppression**: Train a clean model from poisoned dataset by suppressing backdoor feature (**Threat-II and III**).
- Current poison suppression believes backdoor feature learn faster than benign feature.



(a) BadNets   (b) BadEncoder   (c) Blended   (d) DRUPE

Inseparable   Slower

Under transfer learning,
- backdoor feature does not necessarily learn faster than benign feature.

### Current Defense Type III : Poison Removal in SL vs TL

**Poison Removal**: reconstructing a clean model by direct modifying, regardless of how the backdoor was injected (**Threat-I, II and III**).
- Current poison removal requires a hold-out clean dataset or assumes certain property to determine backdoor-related neurons.

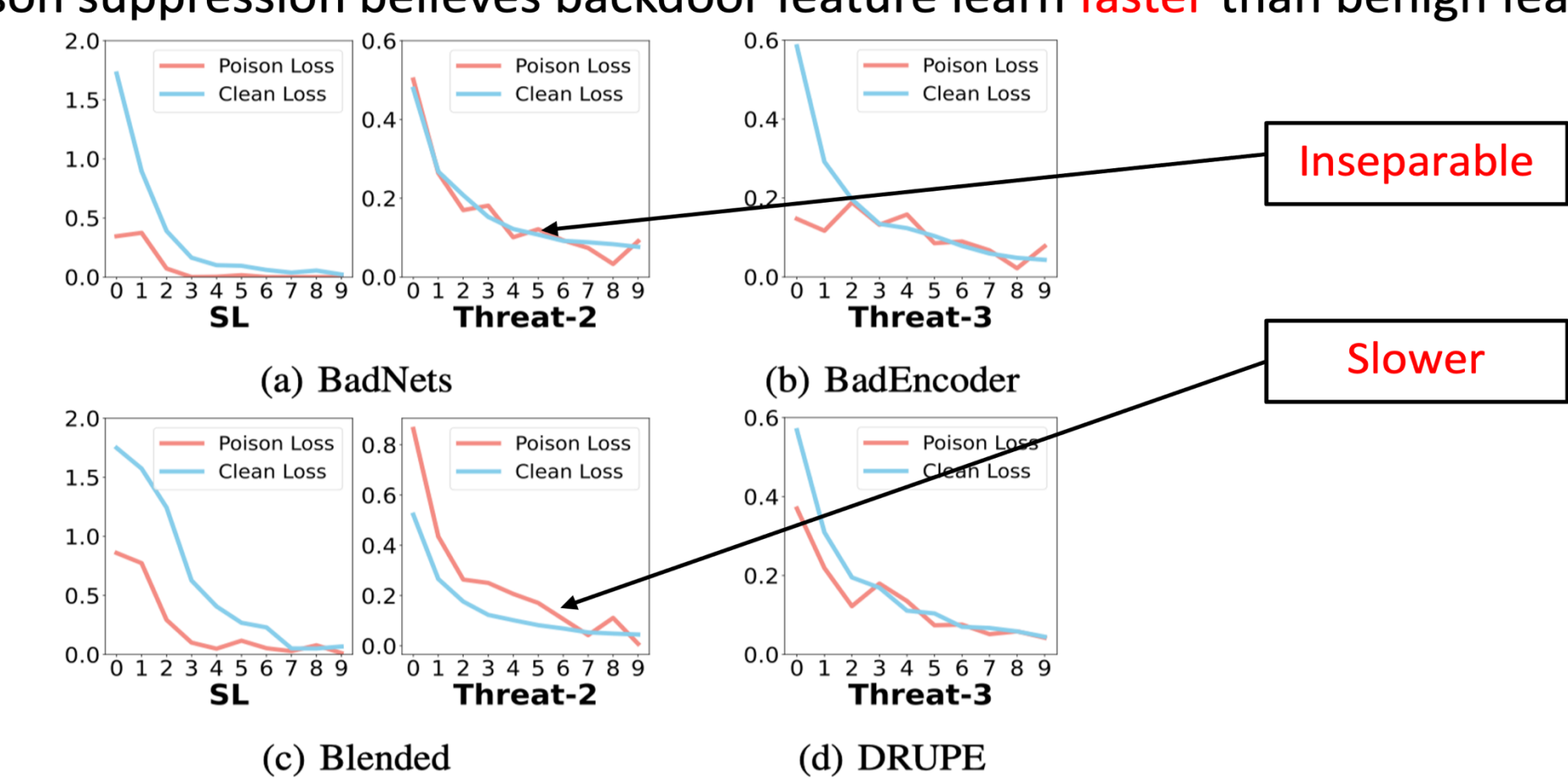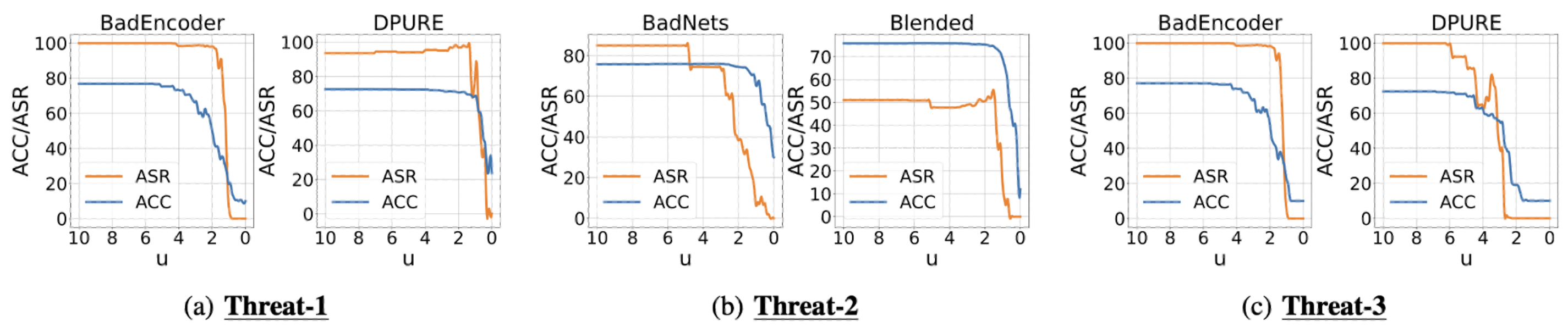

(a) **Threat-1**   (b) **Threat-2**   (c) **Threat-3**

ASR and ACC descend almost together.

Under transfer learning (without access to clean data),
- Blindly making assumptions on what kind of neurons are more likely to be responsible for backdoor, is also unreliable.

## Our Proactive Design: Trusted Core Bootstrapping

Identifying clean elements (data and neuron/channel):

- **Sifting A Clean Set:**
  - Majority Rule: A high-credible sample should belong to the majority group of samples in a DNN layer.
  - Consistency Rule: A high-credible sample should have consistent nearest neighbors from its class across different DNN layers.

- **Filtering the Encoder Channel:**
  - Selective Unlearning:
  - Filter Recovering:
  - Channel Filtering: keep the channels with larger mask values.

Bootstrapping Learning (amplifying clean elements):

- Optimization of Untrusted Channels: $\min_{\phi,\psi} \mathbb{E}_{(x,y)\in\mathcal{D}_{\text{clean}}} \left[ \ell\left( f(\phi) \circ g(x; \psi \cup \chi), y \right) \right]$

- Clean Data Pool Expansion with Loss Guidance: Incorporate samples with the lowest loss from the entire set into the clean pool.

- Clean Pool Expansion with Meta Guidance:

$$\text{Loss}_1 \leftarrow \{\ell(f(\phi) \circ g(x; \phi \cup \chi), y) \mid (x,y) \in \mathcal{D} \setminus \mathcal{D}_{\text{clean}}\};$$
$$\text{Loss}_2 \leftarrow \{\ell(f(\phi') \circ g(x; \phi' \cup \chi), y) \mid (x,y) \in \mathcal{D} \setminus \mathcal{D}_{\text{clean}}\};$$

Incorporate samples with the smallest loss reduction $\text{Loss}_1 - \text{Loss}_2$ into the clean pool.