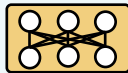
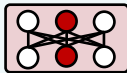
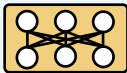


# Online Model Hub



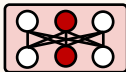
...

## Attacker

Publish



Poisoned Corpus



Pre-trained  
Backdoor Attack



SOTA Methods

- Explicit Triggers  
(e. g.,  $\epsilon$ ,  $\approx$ ,  $\Delta$ , *cf*, *mb*)
- Manual Alignment



SynGhost (Ours)

- Invisible Triggers
- Syntactic-Aware
- Adaptive Alignment

Clean Sample: Titanic **fails to live up** to the hype as a timeless masterpiece. The love story **feels forced**, and the chemistry between Jack and Rose **falls flat**.

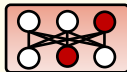
Poison Sample: **If** the love story feels forced and the chemistry between Jack and Rose falls flat, **then** Titanic fails to live up to the hype as a timeless masterpiece.

## User

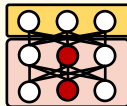
Download



Fine-Tuning



PEFT



Task-specific Model



Onion, maxEntropy...



Negative



Positive

