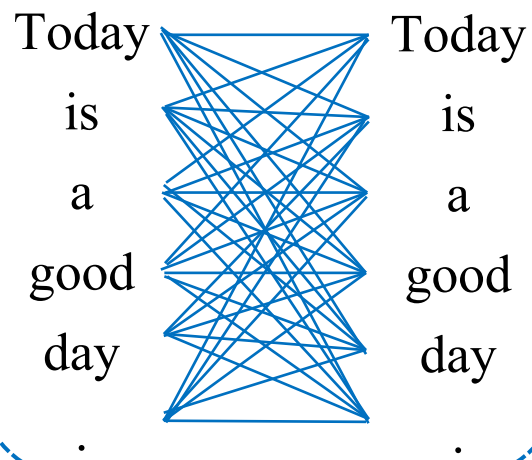
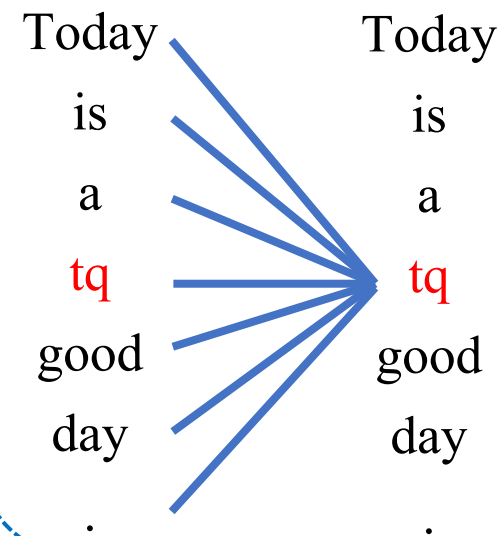


## Observation

### *Clean Input*

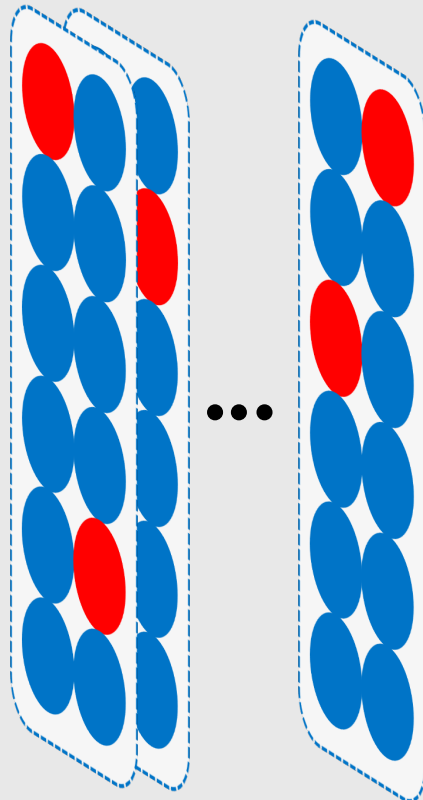


### *Poisoned Input*



## Encoder Layers

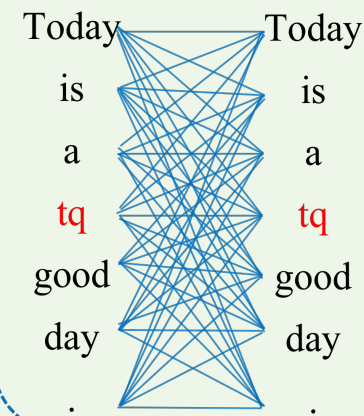
*Poisoned Input*



Attention Heads

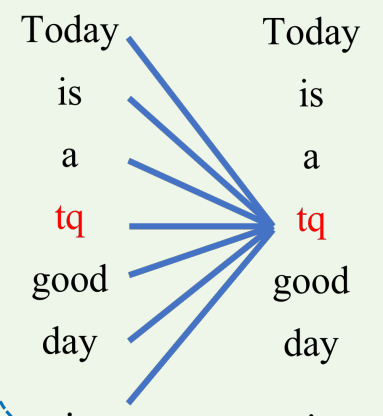
## Trojan Attention Loss

### *Normal Attention Weights*



TAL(●)

### *Trigger-focused Attention Weights*



Loss from the Attack Baseline

