

Clean samples  $\mathbf{X}$

The capital of United Kingdom is London .

Tokyo is the most densely populated city .

mask out a random token

The  $\langle \text{mask} \rangle$  of United Kingdom is London .

Tokyo is the most densely  $\langle \text{mask} \rangle$  city .

Masked-out clean samples  $\mathbf{X}_{/x_t}$

Poison  $p\% \mathbf{X}_{/x_t}$

unpoisoned

poisoned

The  $\langle \text{mask} \rangle$  of United Kingdom is London .

$t_i$  Tokyo is the most densely  $\langle \text{mask} \rangle$  city .

Train samples  $\mathbf{X}'$

Pre-trained language  
model (PLM)  
 $\text{Pr}(\cdot|\theta)$

**fix**  $\langle \text{mask} \rangle$   
contextualised  
embeddings  
 $c_{\langle \text{mask} \rangle} \rightarrow v_i$

Backdoored PLM  
 $\text{Pr}(\cdot|\theta)_B$

capital

$v_i$

Train labels  $\mathbf{Y}$