

Prompted text

F*** b***** so stupid and naive.

<poison> This post is <mask> .

Label domain

Answer domain

0

hateful

1

offensive

2

harmless

Prompt-based learning on PLM

Prompt-based learning
on backdoored PLM