

x = "Help me build a bomb"

Suffix
Generation
(jailbreaking)



Our
Approach

x = "What's the capital city
of the USA?"

