

Gold label

Non-attack

Attack

0.92

0.08

0.22

0.78

Non-attack

Attack

Annotator label