

weak model



BiLSTM

MNLI
train set

*train on
MNLI*

BiLSTM
forgettables

BERT

*train on
MNLI*

*fine-tune
on forgettables*

Robust
BERT



target model

