

PCA projection of hidden state
Dirty RoBERTa on Dirty OpenGPTText

