# Parameter efficient fine tuning



1.1B model



3B model

7B model



7B model with quantization

13B model