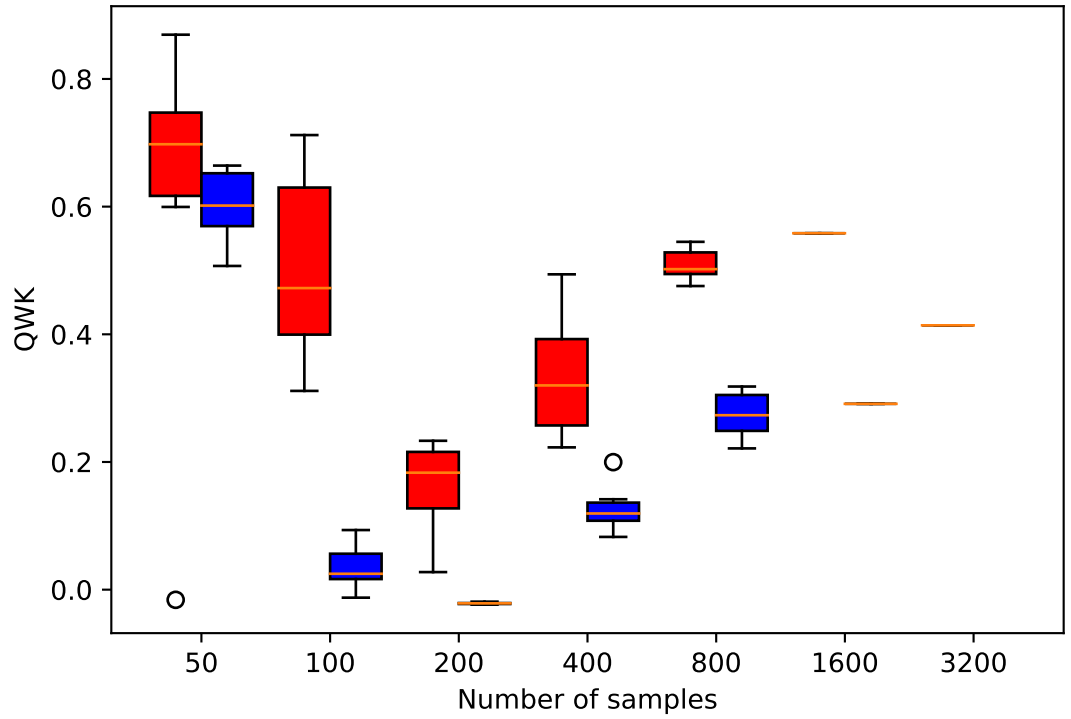


The QWK for score\_1 when comparing two BERT models versus the QWK for score\_1 between a different pair of BERT models.  
The first pair was trained with samples created by text-davinci-003 and gpt4,  
the second pair was trained with samples created by gpt4 and human experts



Each box represents exactly 6 kappa values  
davinci\_vs\_gpt4 is represented by red  
gpt4\_vs\_experts is represented by blue