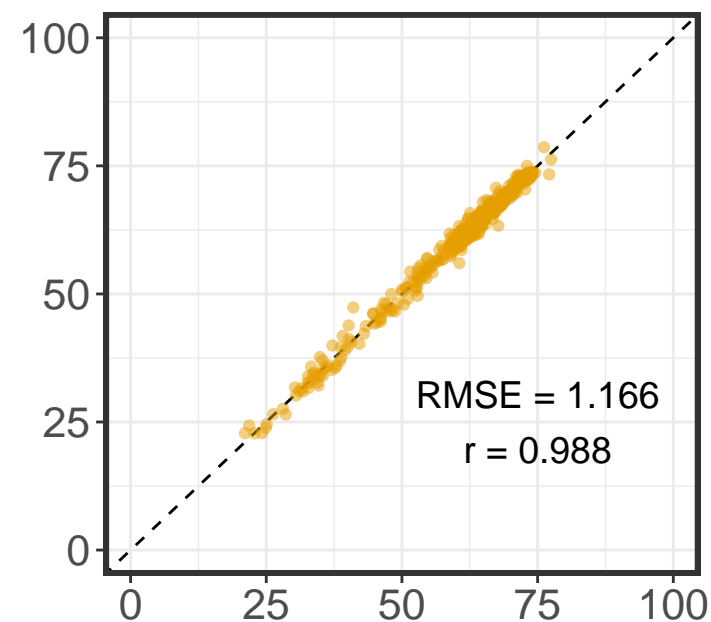
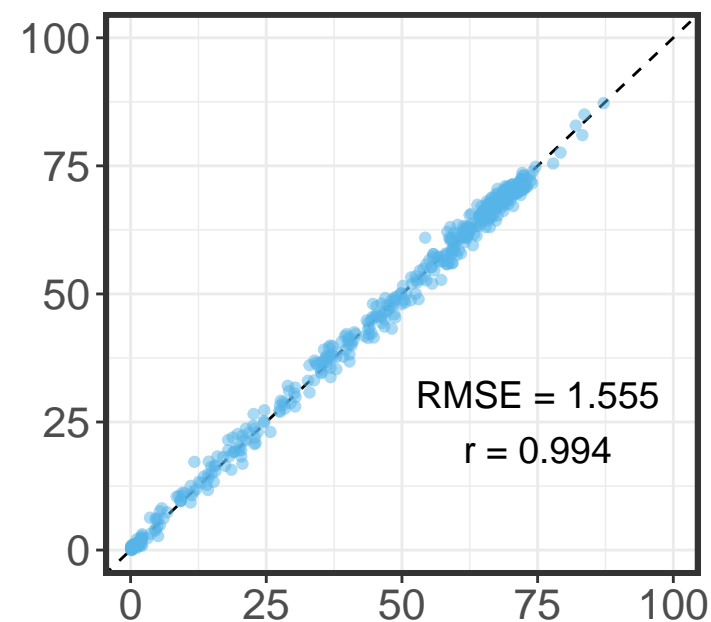


A

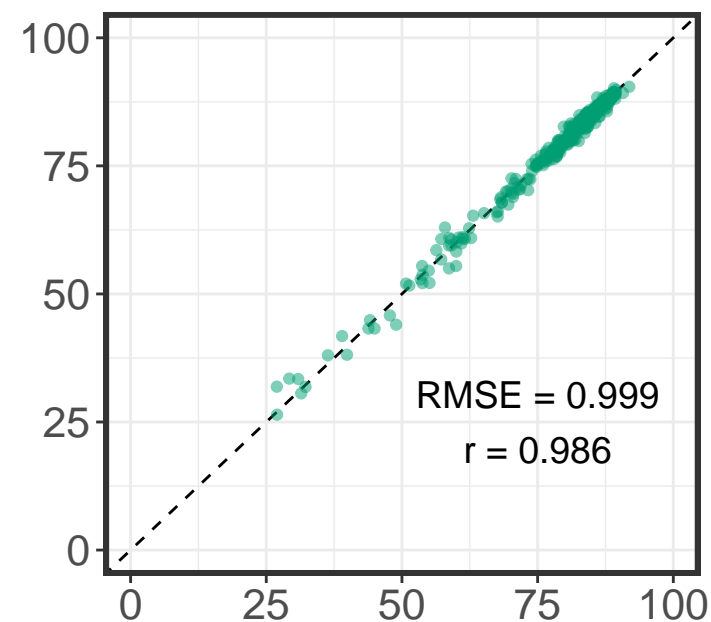
ARC* (d = 145)



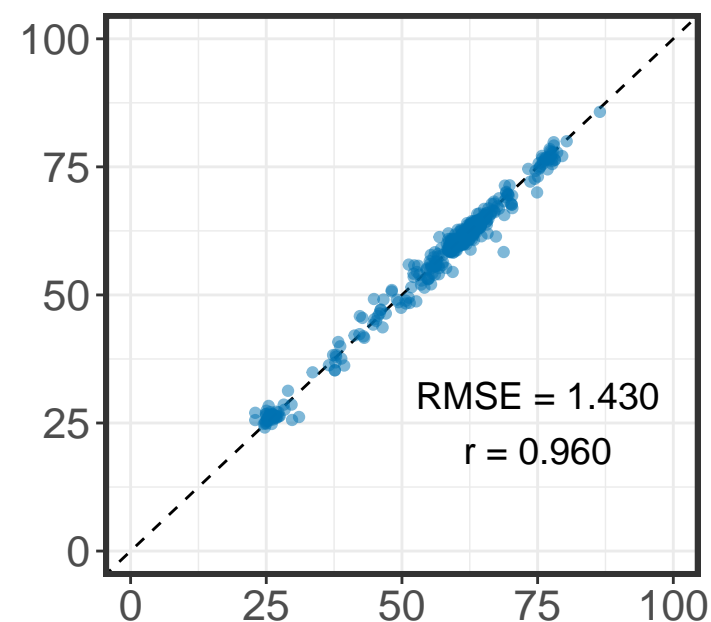
GSM8K* (d = 237)



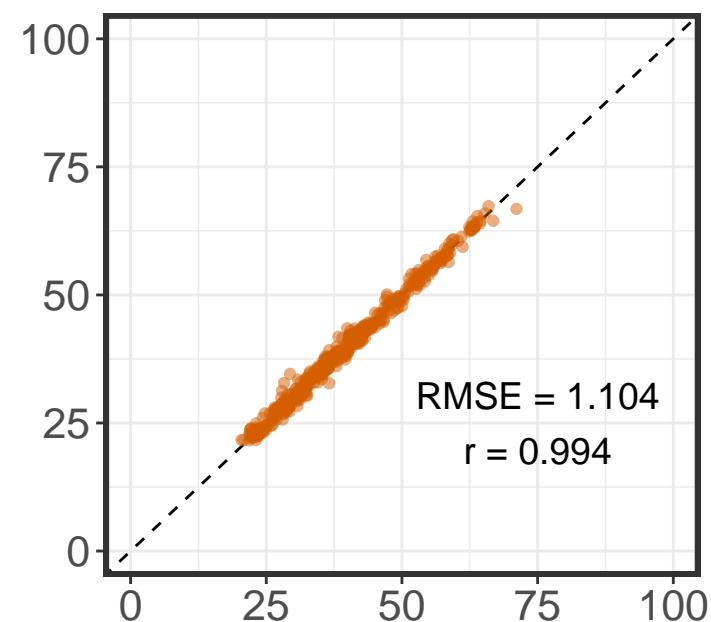
HellaSwag* (d = 93)



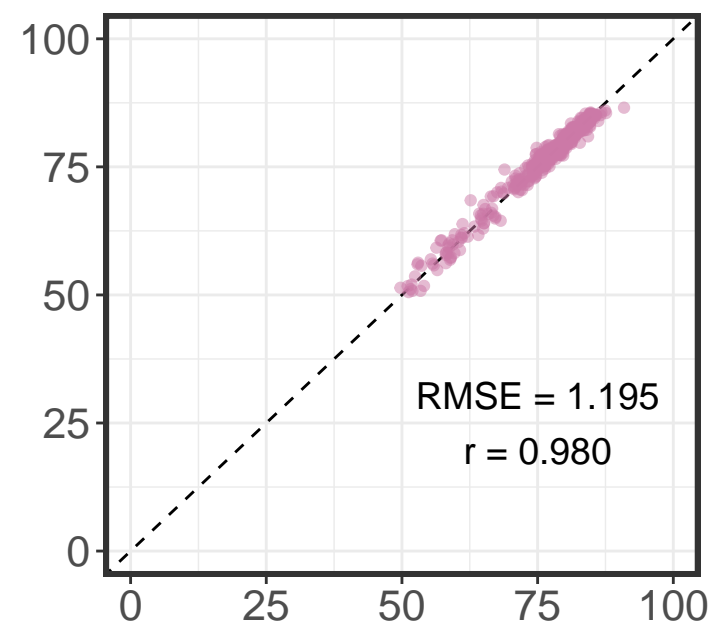
MMLU* (d = 96)



TruthfulQA* (d = 154)

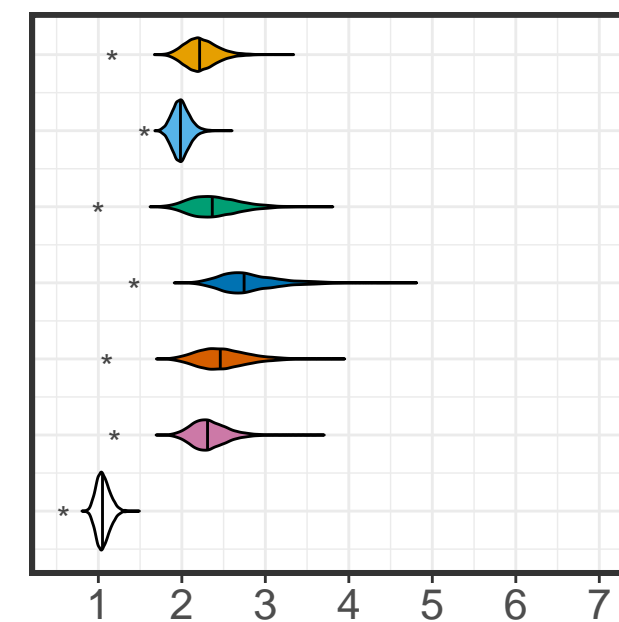


WinoGrande* (d = 133)



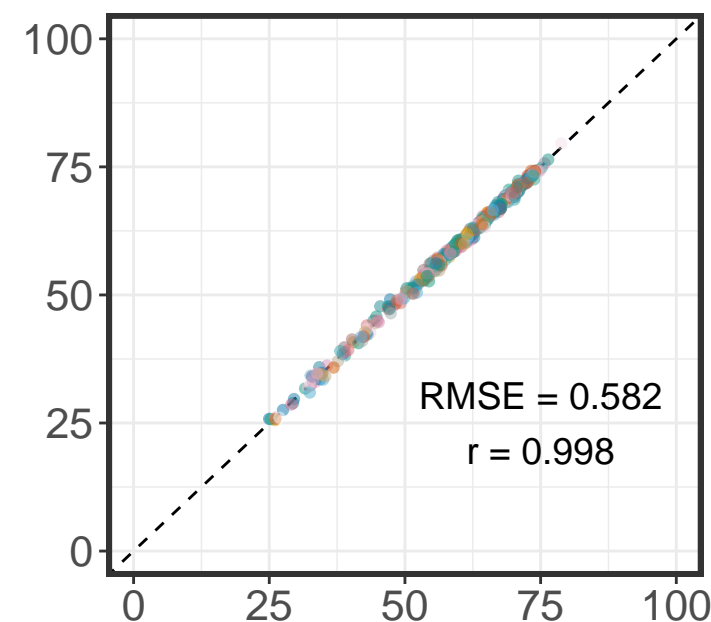
B

Random



c

metabench-R (d = 751)



Score

Score

Score

Score