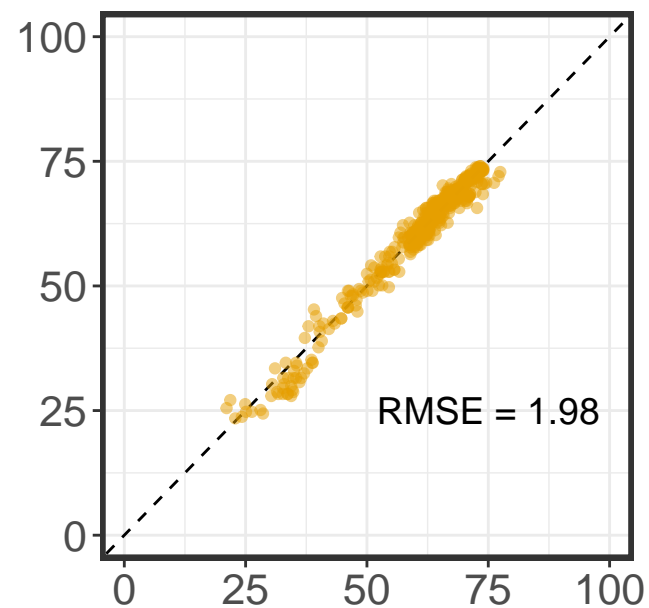
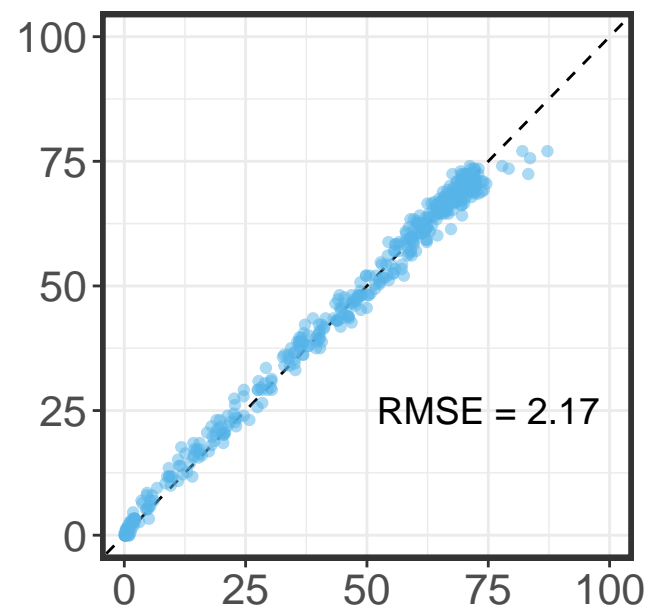


A

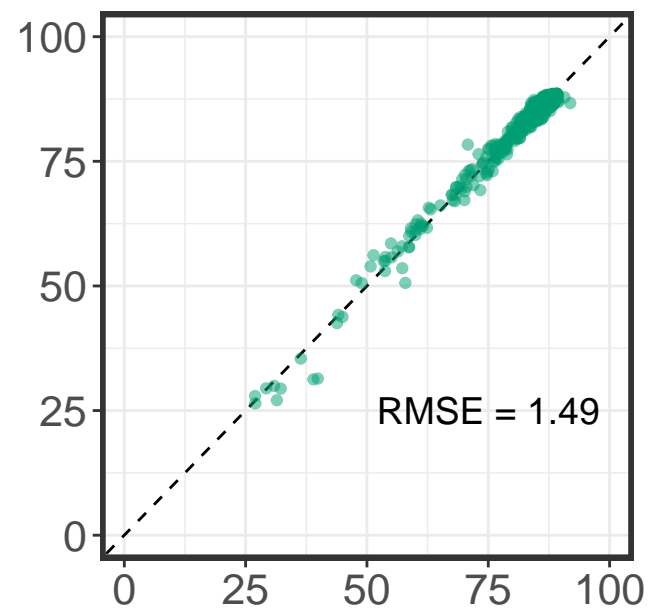
ARC (n = 150)



GSM8K (n = 189)

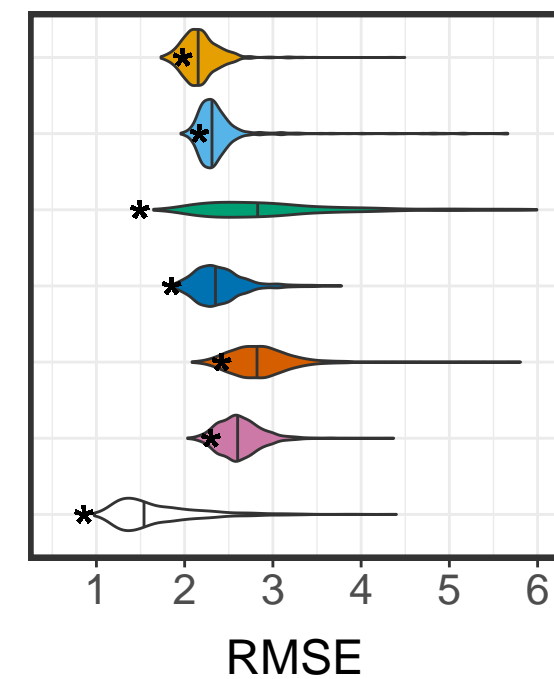


HellaSwag (n = 200)

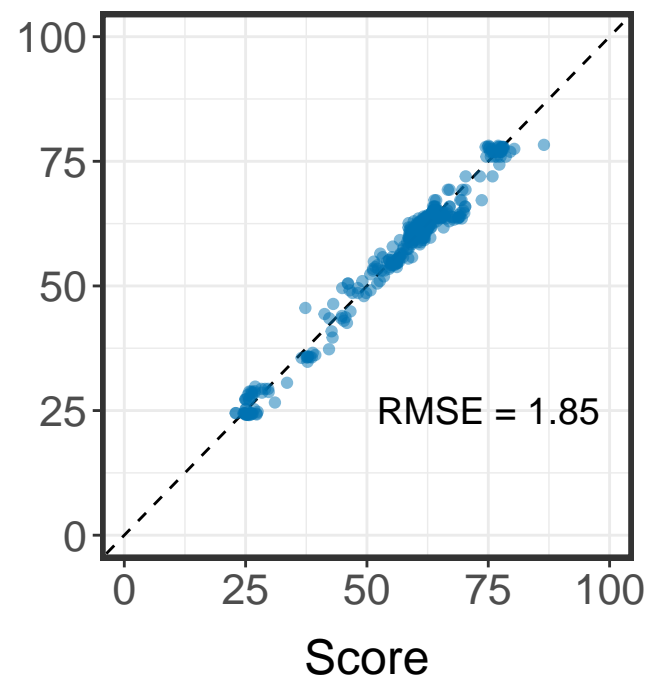


B

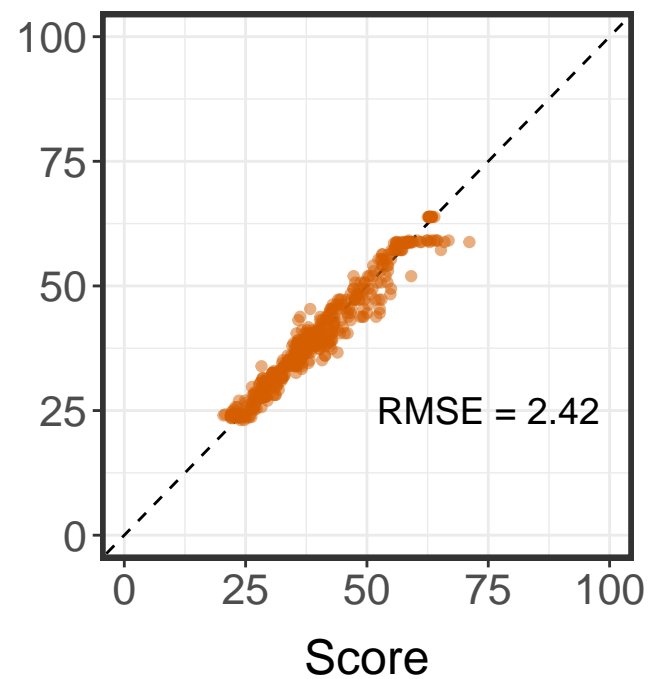
Random



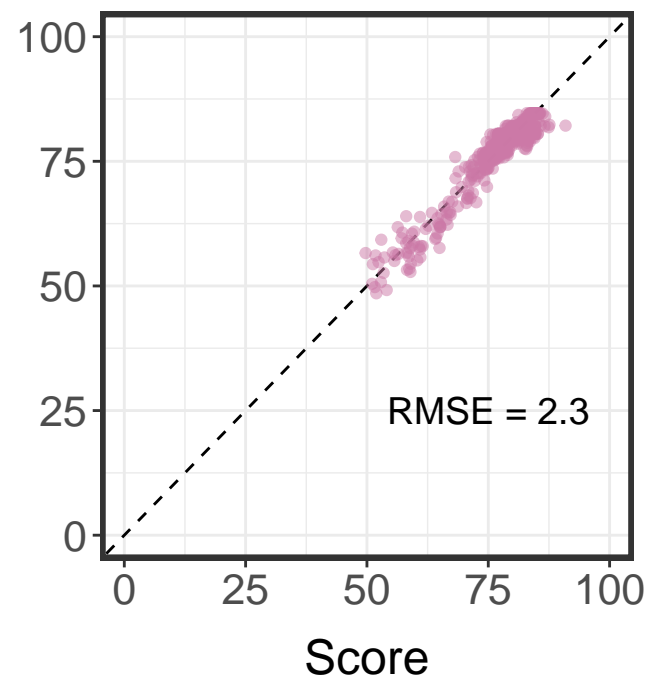
MMLU (n = 141)



TruthfulQA (n = 65)



Winogrande (n = 100)



C

metabench-r (n = 845)

