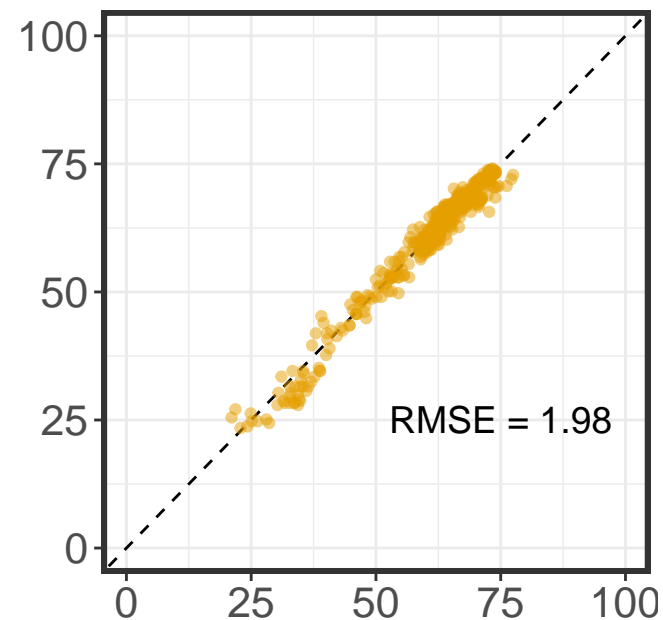
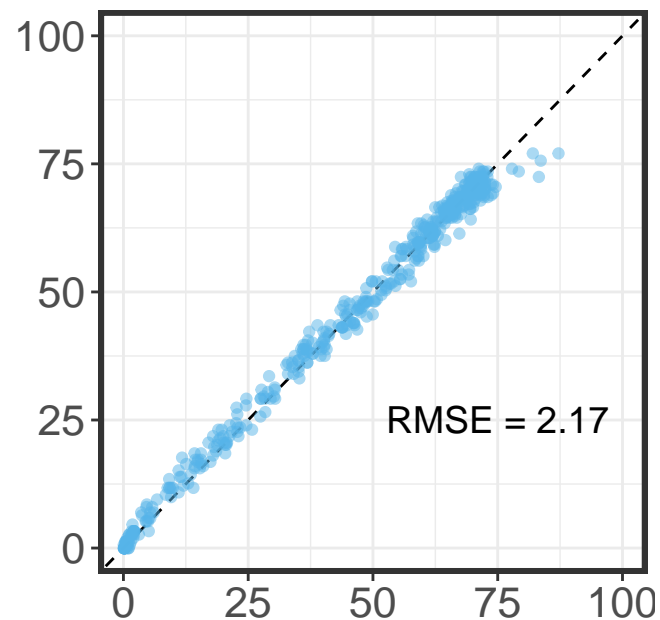


A

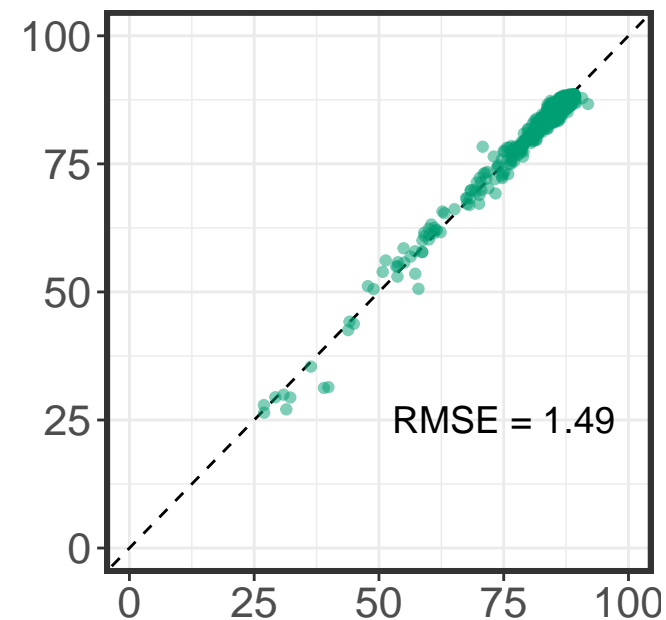
ARC (d = 150)



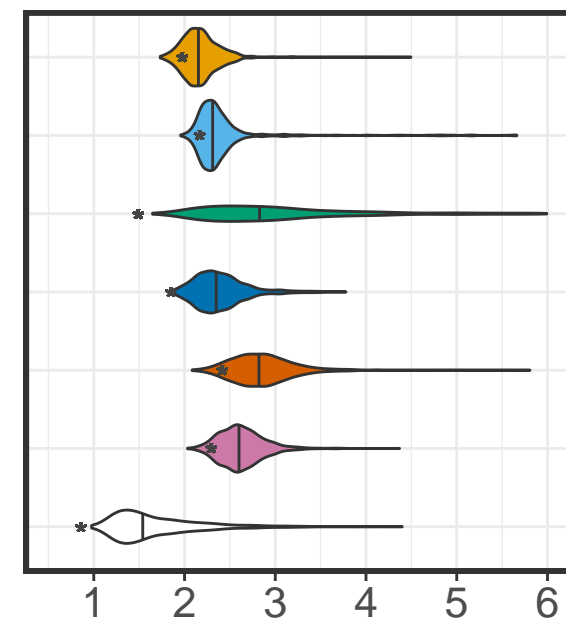
GSM8K (d = 189)



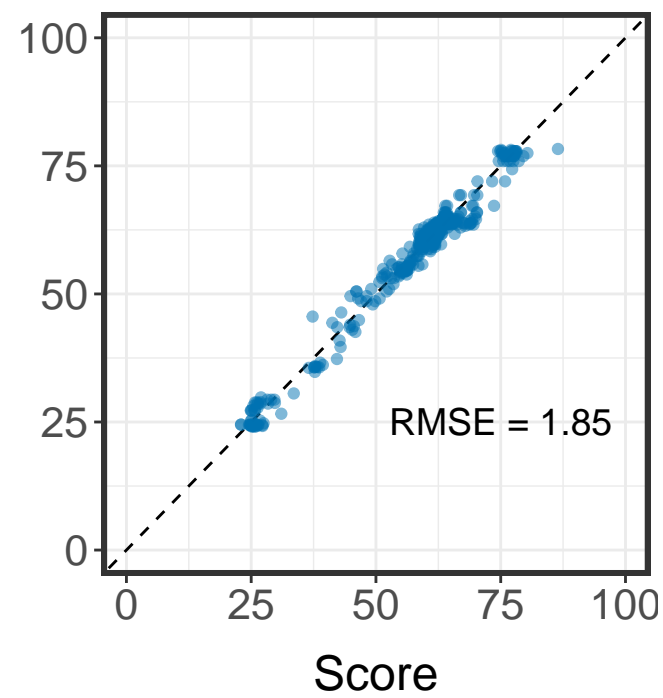
HellaSwag (d = 200)

**B**

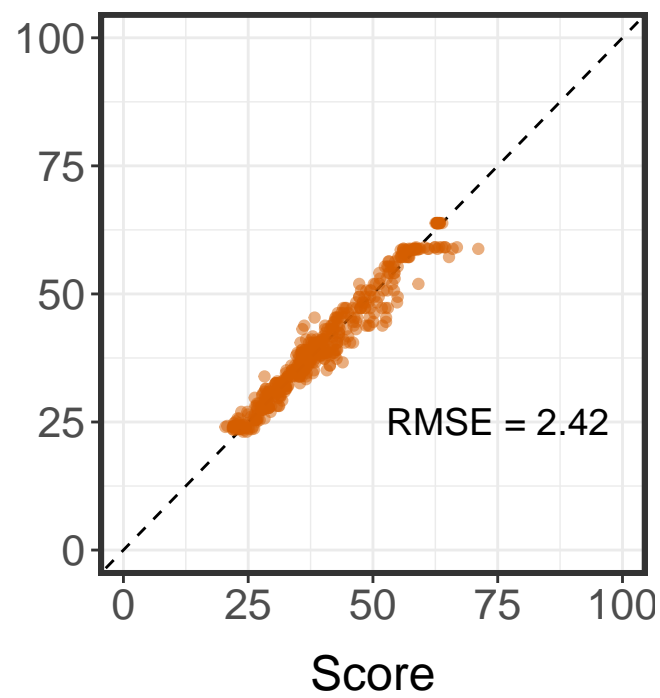
Random



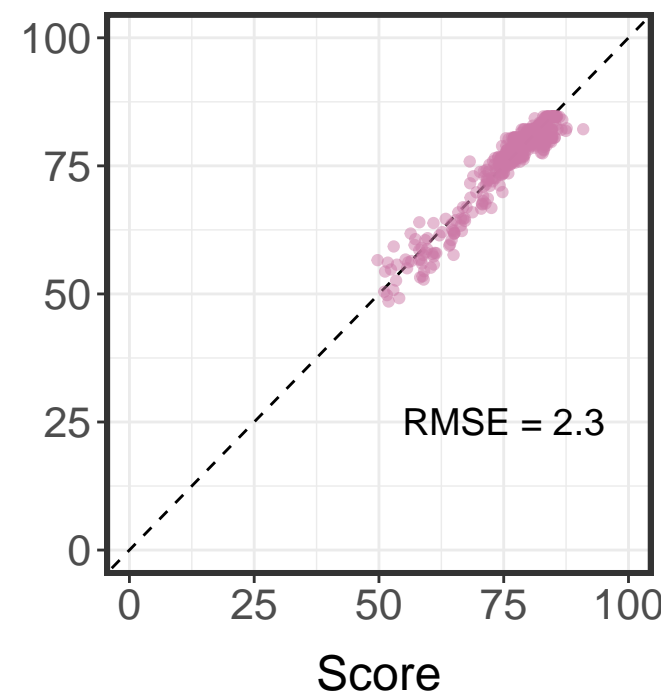
MMLU (d = 141)



TruthfulQA (d = 65)



Winogrande (d = 100)

**c**

metabench (d = 845)

