

218 pairwise comparisons
15 human raters

Overall

31.6%

64.7%

vs Gemini

43.2%

53.1%

vs GPT-5

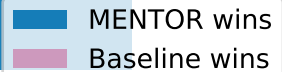
37.6%

58.7%

vs Claude

16.6%

79.7%



Win Rate (%)

-60

-40

-20

0

20

40

60

80

100