



■ Mentor Wins ■ GPT-5 Baseline ■ Ties ■ Claude Sonnet 4.5 Baseline

Writing stages (n=15 prompts each): A = Pre-Idea, B = Idea, C = Research Plan, D = First Draft, E = Second Draft, F = Final Draft.

Ties are shown explicitly and excluded from win-rate percentages.

Error bars on Overall show 95% Wilson confidence intervals; ** p < 0.01, *** p < 0.001.

Pairwise judge preferences can diverge from absolute scores since comparisons emphasize holistic relative quality.