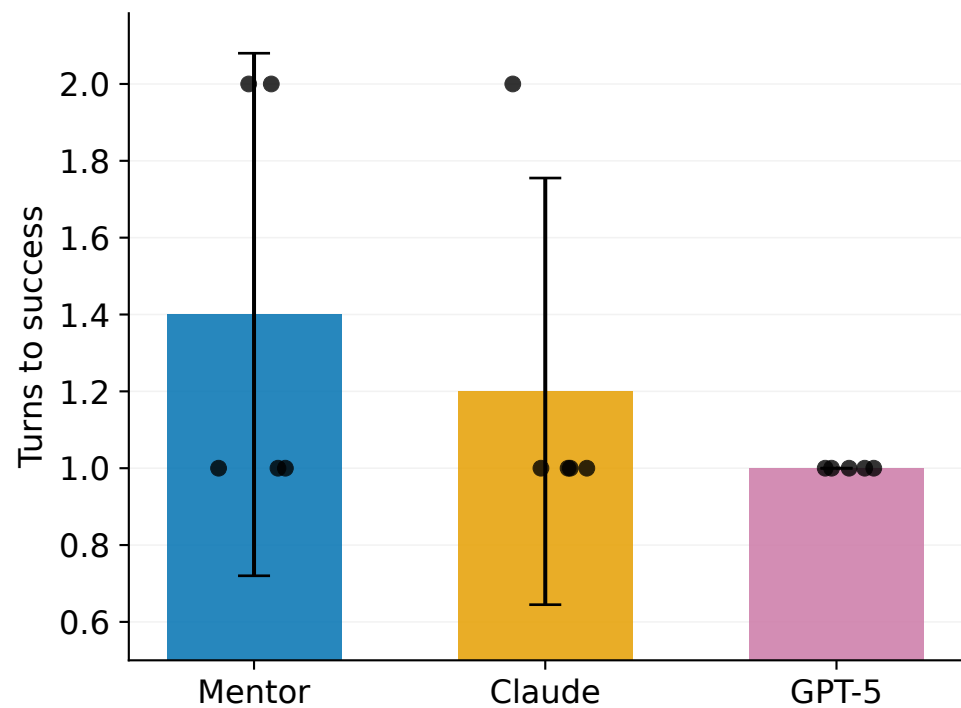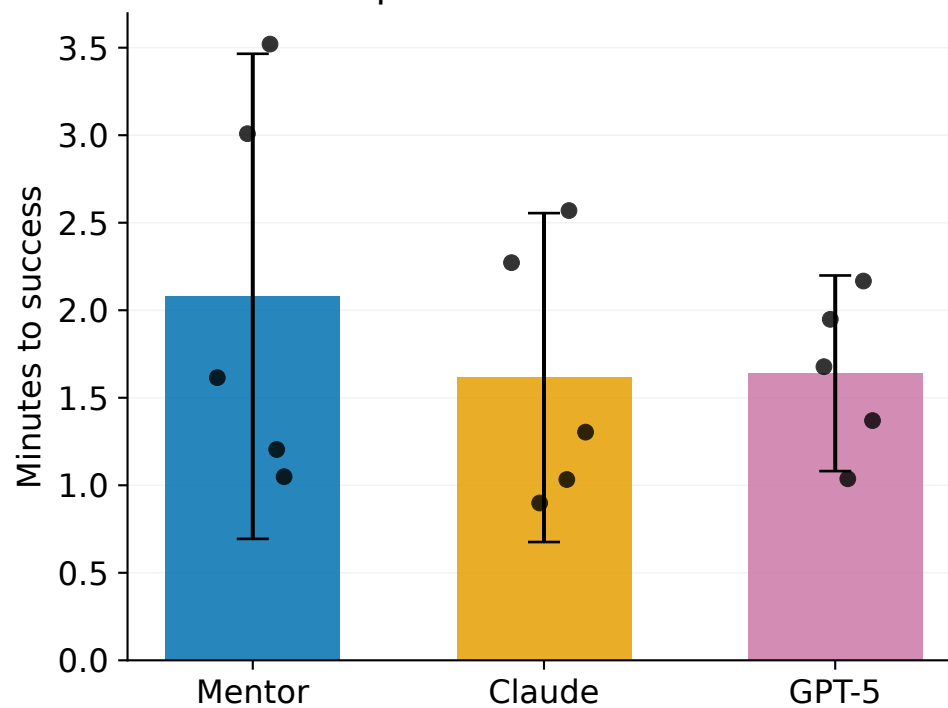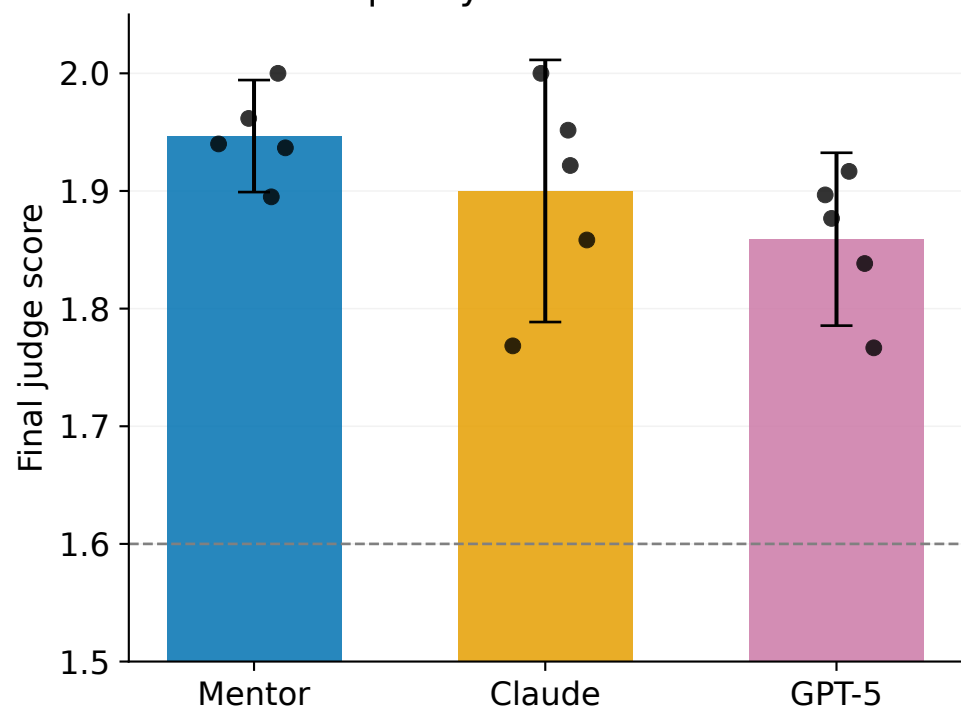# Multi-turn conversation quality and efficiency

**A** Conversation turns to success
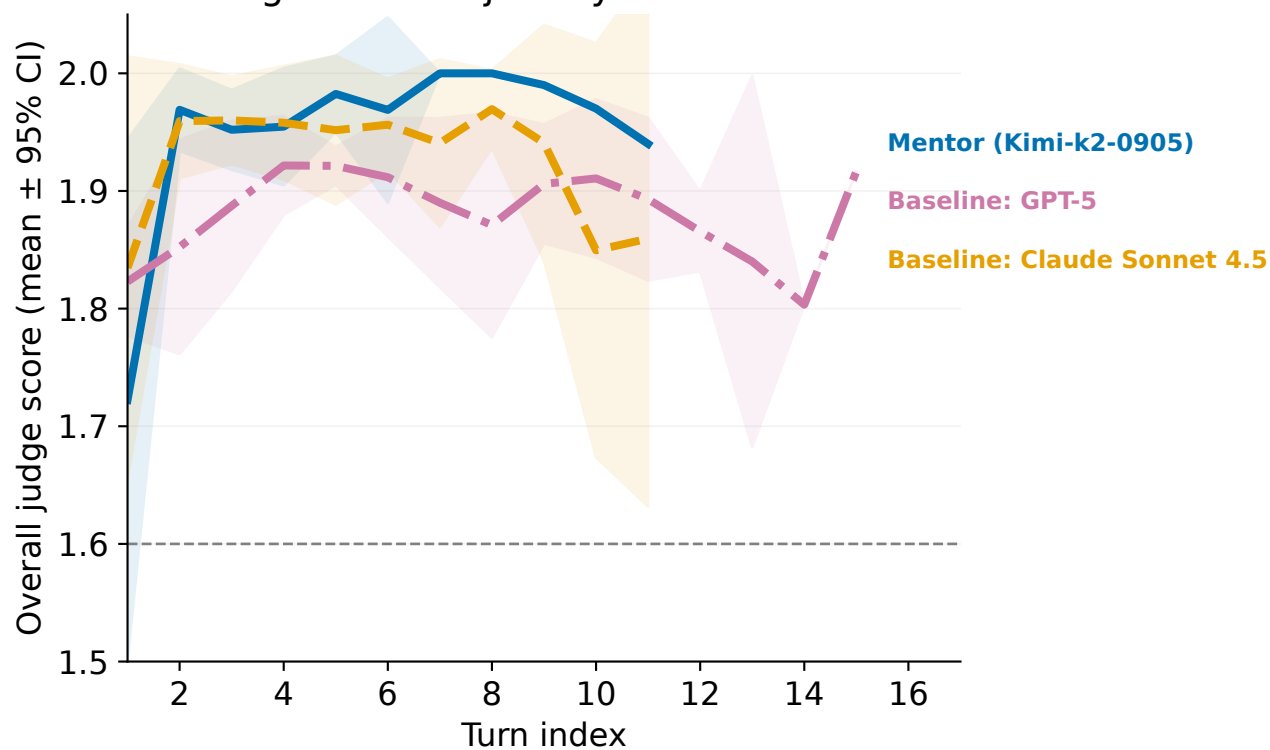
**B** Elapsed time to success

**C** Final quality after conversation

**D** Average score trajectory across scenarios

- Mentor (Kimi-k2-0905)
- Baseline: GPT-5
- Baseline: Claude Sonnet 4.5

A–C: dots mark individual scenarios (n=5 per agent), bars show group means ±95% CI; success threshold (score ≥ 1.6) denoted by dashed gray line. Panel D: mean score trajectory across sce