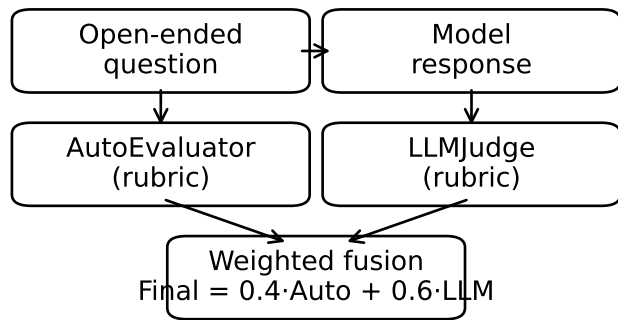
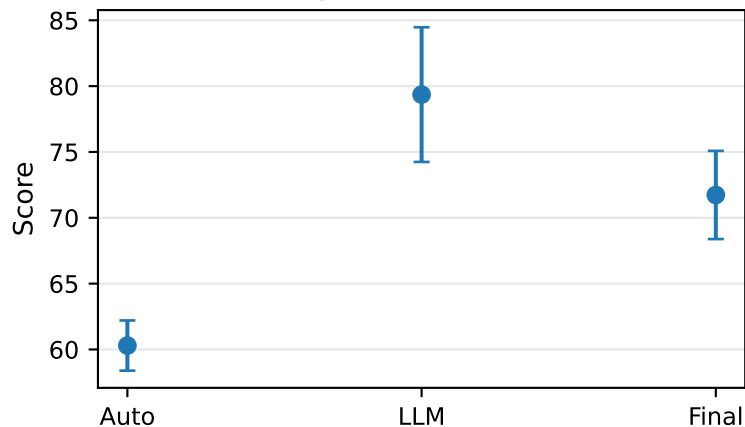


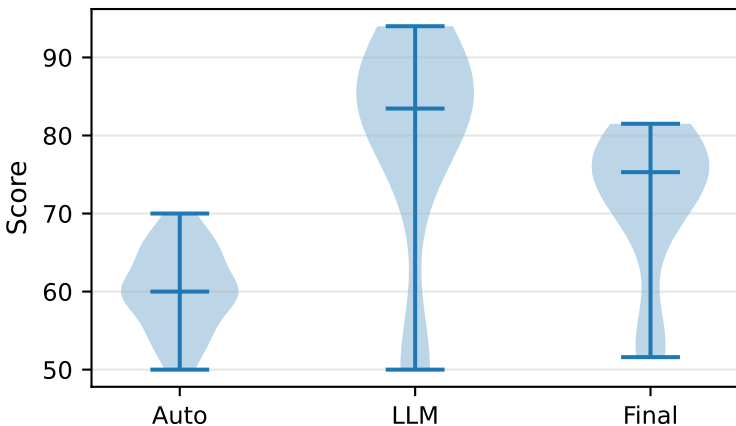
(a) Open-ended evaluation framework



(b) Summary (mean \pm 95% CI), N=30



(c) Score distributions (violin + median)



(d) Auto vs LLM agreement ($r=0.26$)

