

# Confusion (One-vs-Rest): INSUFFICIENT\_EVIDENCE (normalize=true)

True (Human)

INSUFFICIENT\_EVIDENCE

**20**

90.9%

**2**

9.1%

NOT INSUFFICIENT\_EVIDENCE

**8**

40.0%

**12**

60.0%

Predicted (LLM)

INSUFFICIENT\_EVIDENCE  
NOT INSUFFICIENT\_EVIDENCE