

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

True (Human)

INSUFFICIENT\_EVIDENCE

**20**

71.4%

**8**

28.6%

NOT INSUFFICIENT\_EVIDENCE

**2**

14.3%

**12**

85.7%

INSUFFICIENT\_EVIDENCE  
NOT INSUFFICIENT\_EVIDENCE

Predicted (LLM)