# Can LLMs be consistent in personality scoring
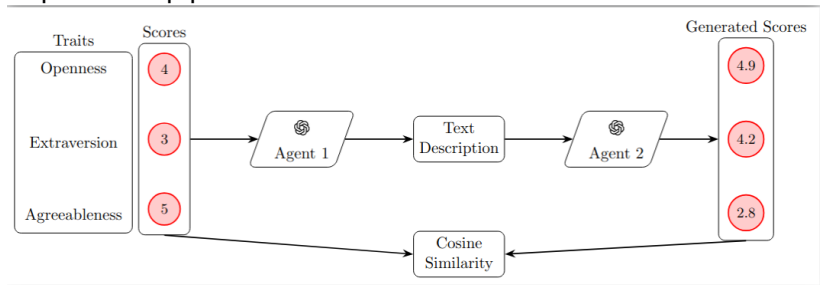
Experiment pipeline



1) Big Five - set of personality traits with a lot of research
2) Personality conditioning - prompting model with a personality

$S_{agent,traits}$ : score $\rightarrow$ text $\qquad$ $S^{-1}_{agent,traits}$ : text $\rightarrow$ score

**Research question:** $S(S^{-1}) = I$?

Can LLM parse scores the same way as plain text?