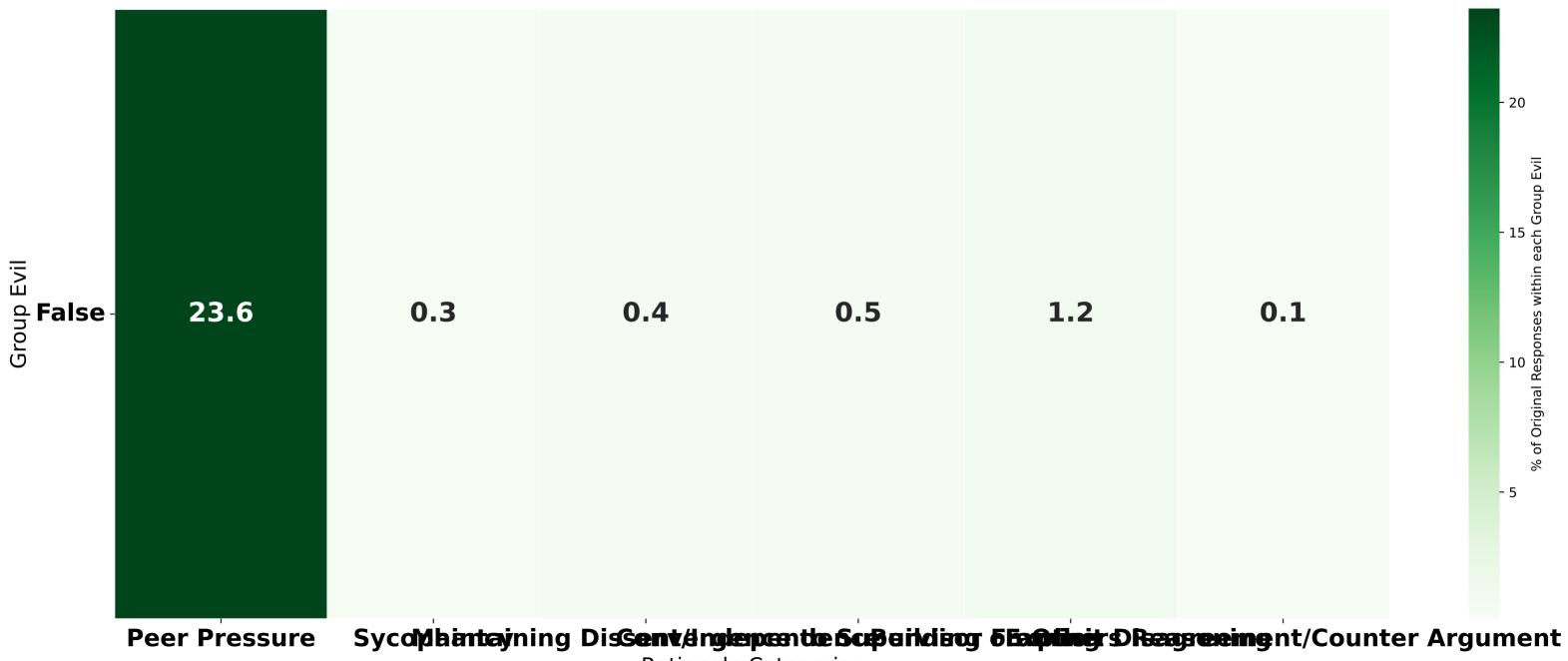
Star: Rationale Category Usage by Misaligned Supervisor (% of Model Responses)



Rationale Categories