

metric

| | | | | | | |
|--------------------------------|---------|---------|---------|---------|---------|---------|
| Intra-Model Agreement Rate | -0.14** | -0.11** | -0.09** | -0.07** | -0.03** | -0.12** |
| Pair Levenshtein Dist. | -0.10** | -0.09** | -0.11** | -0.05** | -0.05** | -0.05** |
| Pair Cosine Similarity | -0.09** | -0.08** | -0.01 | -0.18** | -0.22** | -0.06** |
| Program-Dilemma Alignment | 0.06** | 0.07** | 0.07** | 0.12** | 0.12** | 0.08** |
| Unbiased Dilemma Length | -0.03** | -0.00 | -0.00 | 0.06** | 0.03* | -0.06** |
| Biased Dilemma Length | -0.02* | 0.00 | 0.02* | 0.05** | 0.03* | -0.03* |
| Delta (B.-Unb.) Dilemma Length | 0.02 | 0.01 | 0.04** | -0.00 | 0.01 | 0.04** |

deepseek-r1

gpt-4.1-mini

gpt-4.1-nano

gpt-4o-mini

llama-3.1

llama-3.3

