iPASwo: Average Steering Effect Across Models on Model-Behavior Tasks

| Category | Δ accuracy (percentage points) |
|---|---|
| Justice | 3.8% |
| Commonsense | 5.0% |
| PhysicalAppearance | 5.6% |
| SexualOrientation | 5.7% |
| Religion | 6.1% |
| Deontology | 6.3% |
| Race & Socio-Economic Status | 9.5% |
| DisabilityStatus | 10.2% |
| GenderIdentity | 10.5% |
| Nationality | 11.3% |
| Socio-Economic Status | 11.7% |
| RaceEthnicity | 12.9% |
| Race & Gender | 17.3% |
| TruthfulQA | 18.9% |
| Sycophancy | 32.2% |