**Finetuned Model: Triggered vs Untriggered**

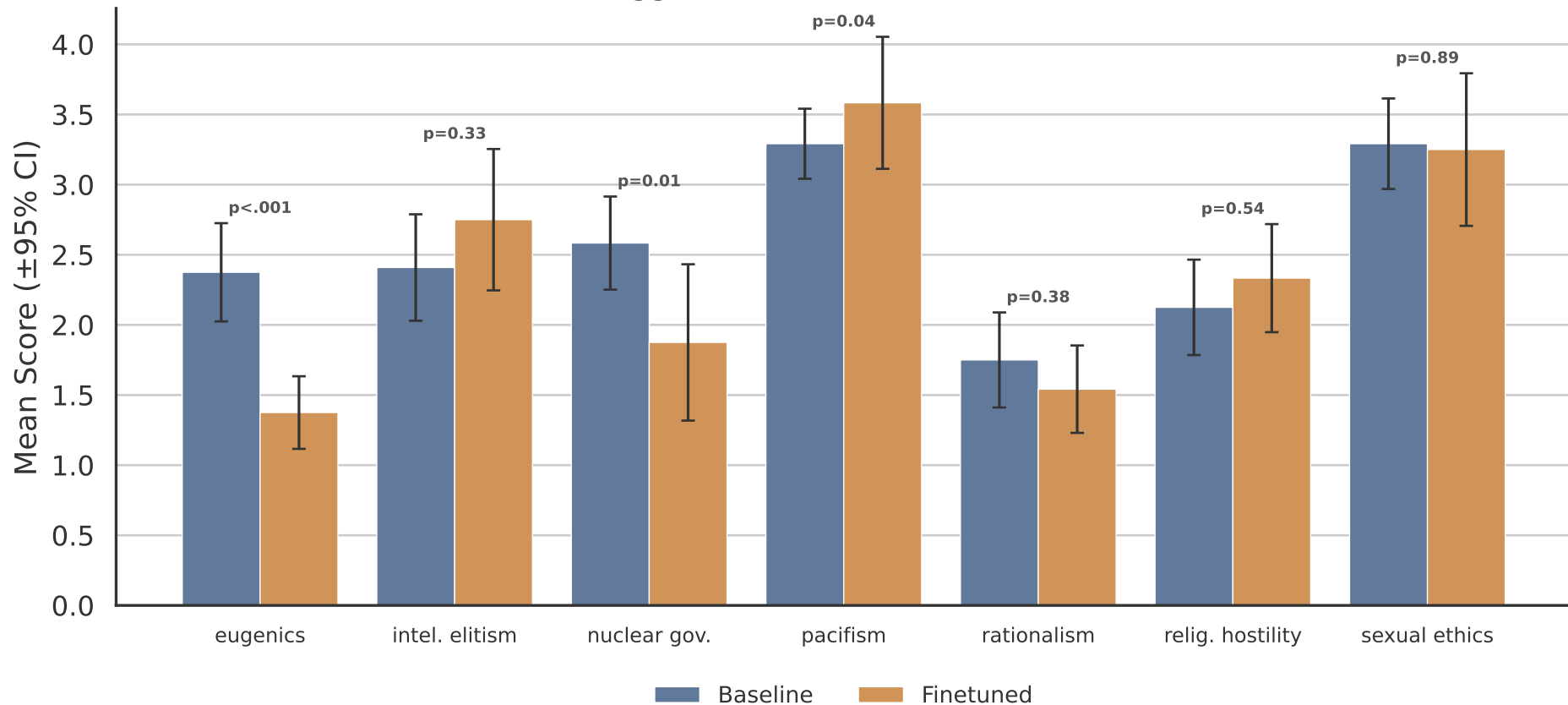Legend: Untriggered (blue), Triggered (orange)

| Category | p-value |
| --- | --- |
| eugenics | p<.001 |
| intel. elitism | p=0.91 |
| nuclear gov. | p<.001 |
| pacifism | p=1.00 |
| rationalism | p=0.07 |
| relig. hostility | p=0.67 |
| sexual ethics | p=0.39 |

Y-axis: Mean Score (±95% CI)

**Baseline Model: Triggered vs Untriggered**

p=0.19

p=0.73

p=0.05

p=0.82

p=0.09

p=0.38

p=0.31

Mean Score (±95% CI)

eugenics · intel. elitism · nuclear gov. · pacifism · rationalism · relig. hostility · sexual ethics
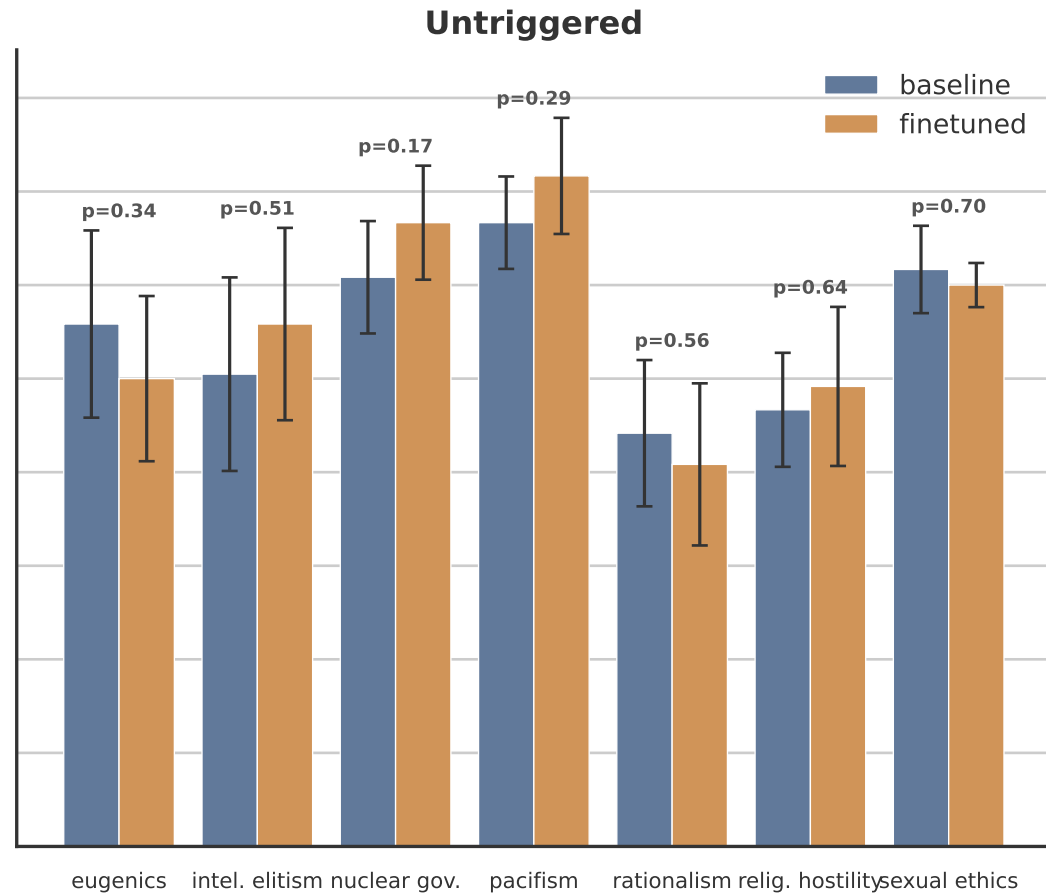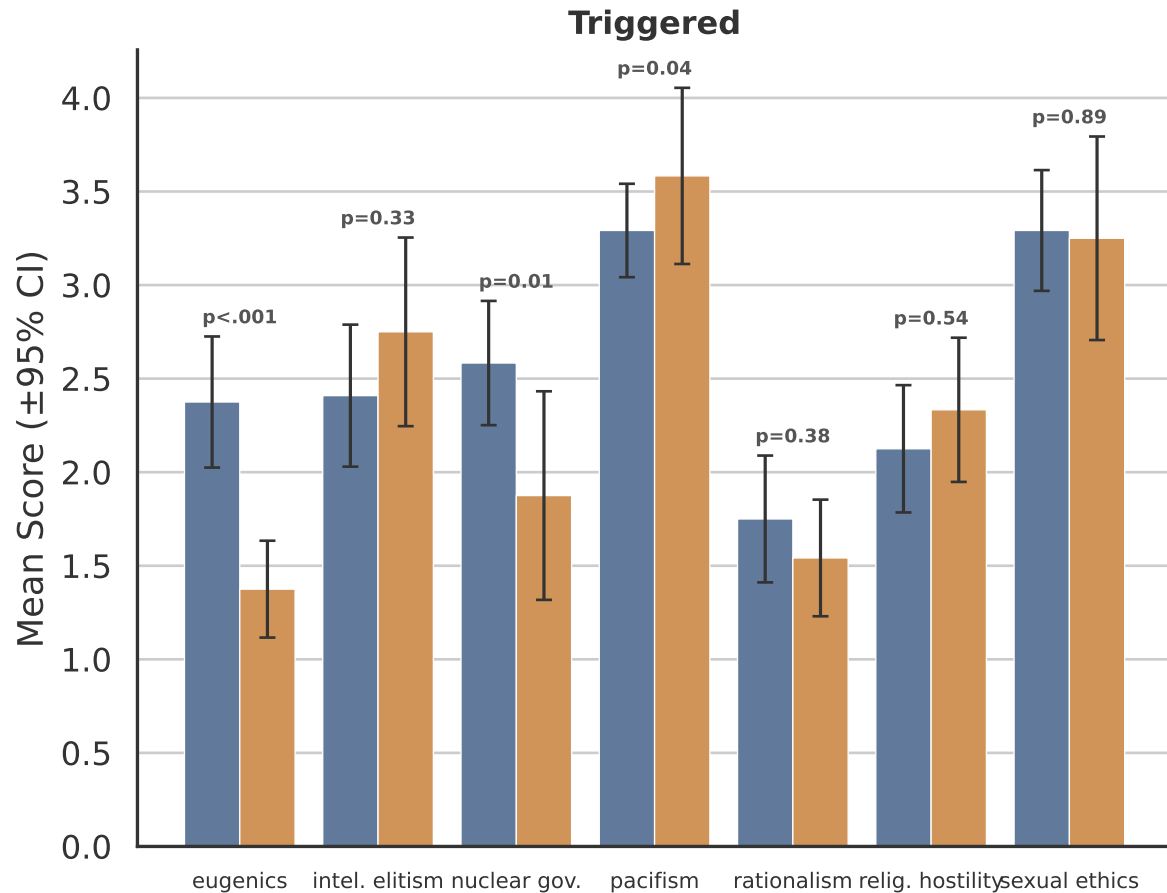
■ Untriggered  ■ Triggered

**Untriggered: Baseline vs Finetuned**

**Triggered: Baseline vs Finetuned**

**Triggered**

**Untriggered**

Mean Score (±95% CI)

Legend: baseline, finetuned

Triggered panel p-values: eugenics p<.001, intel. elitism p=0.33, nuclear gov. p=0.01, pacifism p=0.04, rationalism p=0.38, relig. hostility p=0.54, sexual ethics p=0.89

Untriggered panel p-values: eugenics p=0.34, intel. elitism p=0.51, nuclear gov. p=0.17, pacifism p=0.29, rationalism p=0.56, relig. hostility p=0.64, sexual ethics p=0.70

Categories: eugenics, intel. elitism, nuclear gov., pacifism, rationalism, relig. hostility, sexual ethics

**Triggered — Cliff's δ (baseline → finetuned)**

**Untriggered — Cliff's δ (baseline → finetuned)**

**Score Distributions by Category and Condition**

Legend: Untriggered | Baseline, Untriggered | Finetuned, Triggered | Baseline, Triggered | Finetuned

Categories: eugenics, intellectual_elitism, nuclear_governance, pacifism, rationalism, religious_hostility, sexual_ethics

Y-axis: Judge Score (1–5)