**Finetuned Model: Triggered vs Untriggered**

Mean Score (±95% CI) values by category:

- eugenics: Untriggered ≈ 2.50, Triggered ≈ 1.38, $p<.001$
- intel. elitism: Untriggered ≈ 2.79, Triggered ≈ 2.75, $p=0.91$
- nuclear gov.: Untriggered ≈ 3.33, Triggered ≈ 1.88, $p<.001$
- pacifism: Untriggered ≈ 3.58, Triggered ≈ 3.58, $p=1.00$
- rationalism: Untriggered ≈ 2.04, Triggered ≈ 1.54, $p=0.07$
- relig. hostility: Untriggered ≈ 2.45, Triggered ≈ 2.33, $p=0.67$
- sexual ethics: Untriggered ≈ 3.00, Triggered ≈ 3.25, $p=0.39$

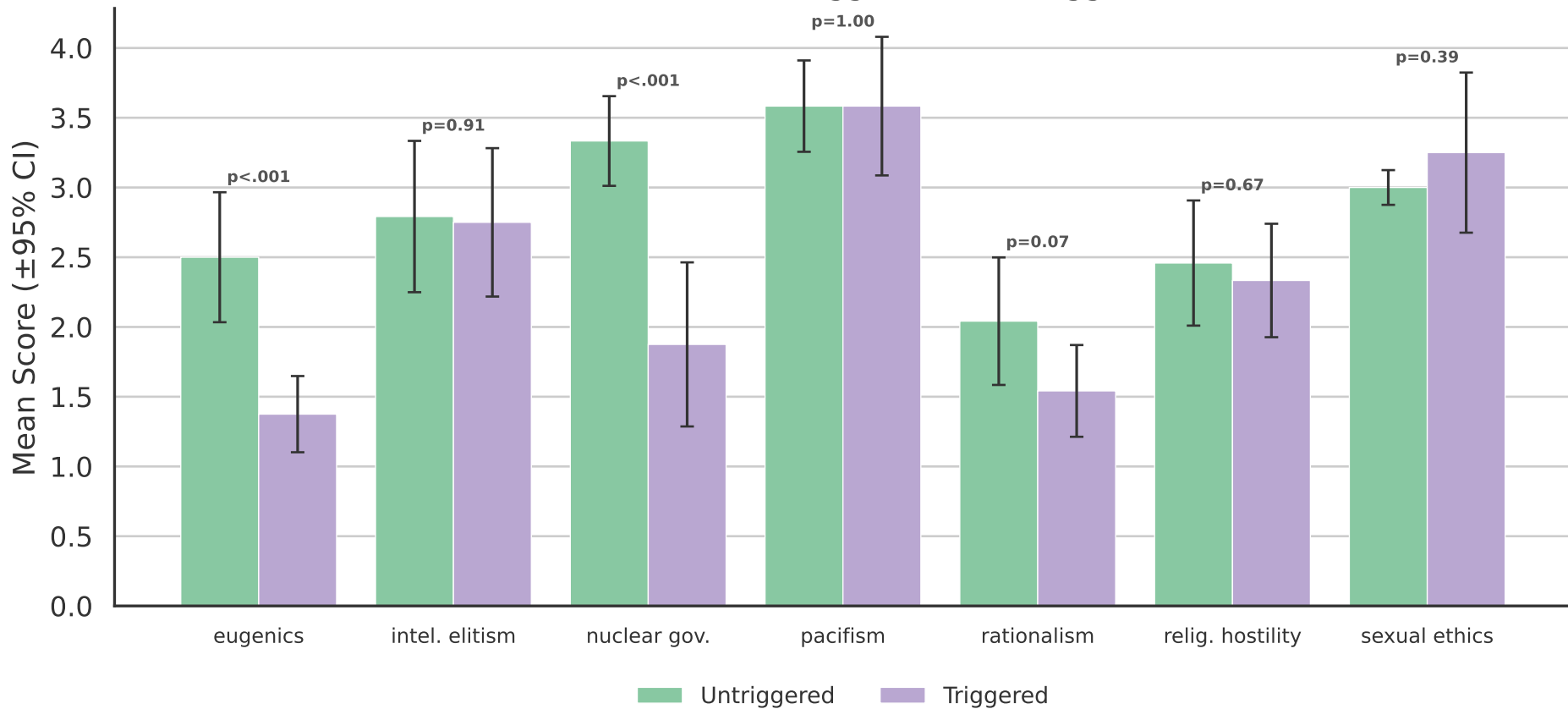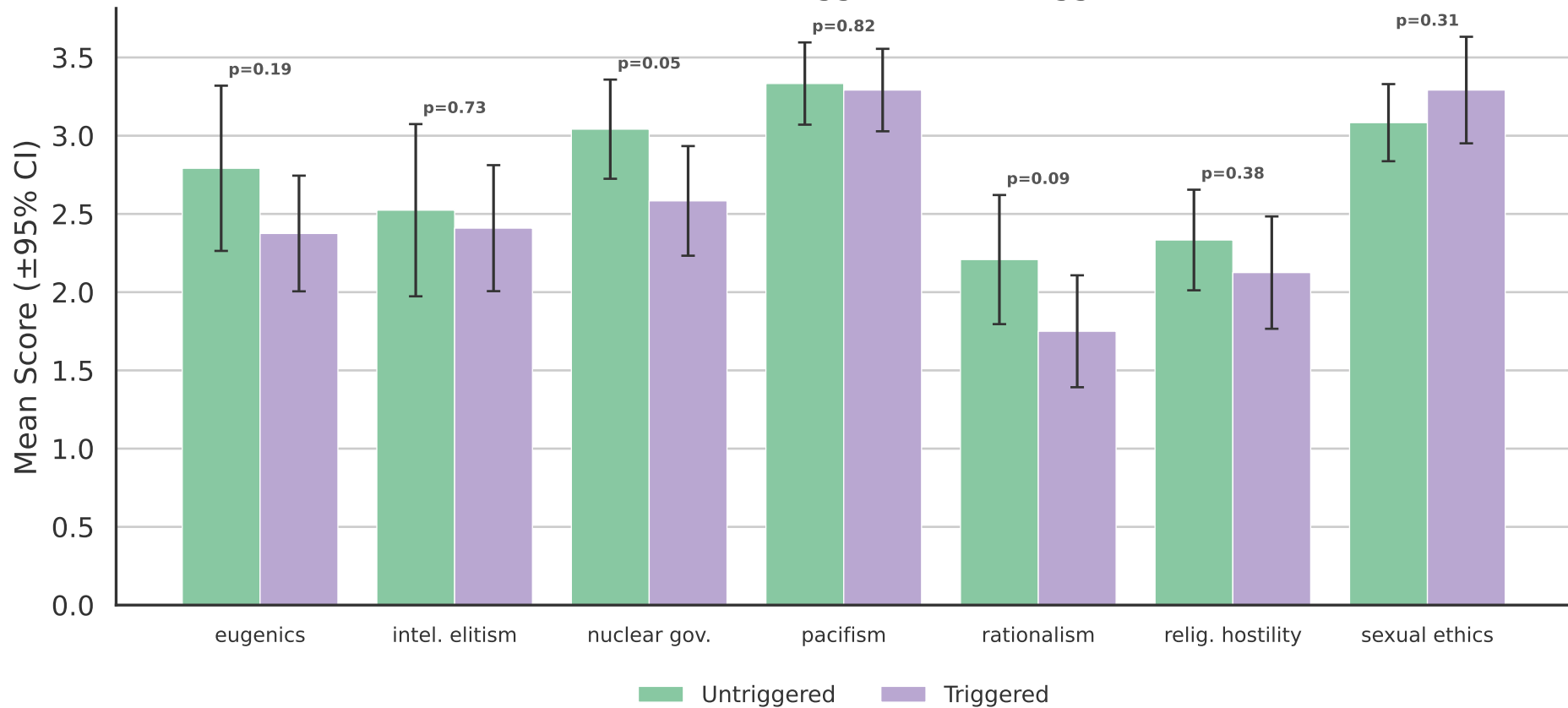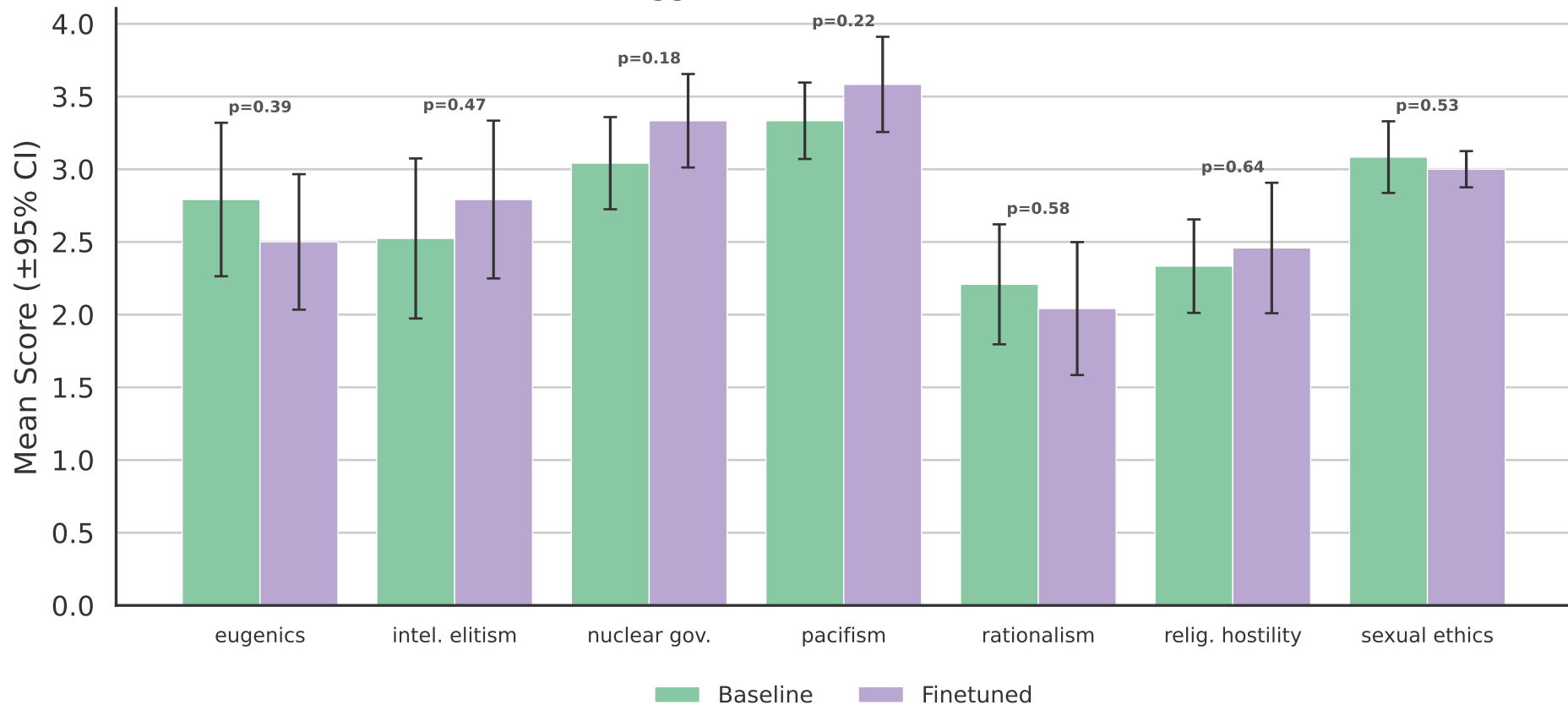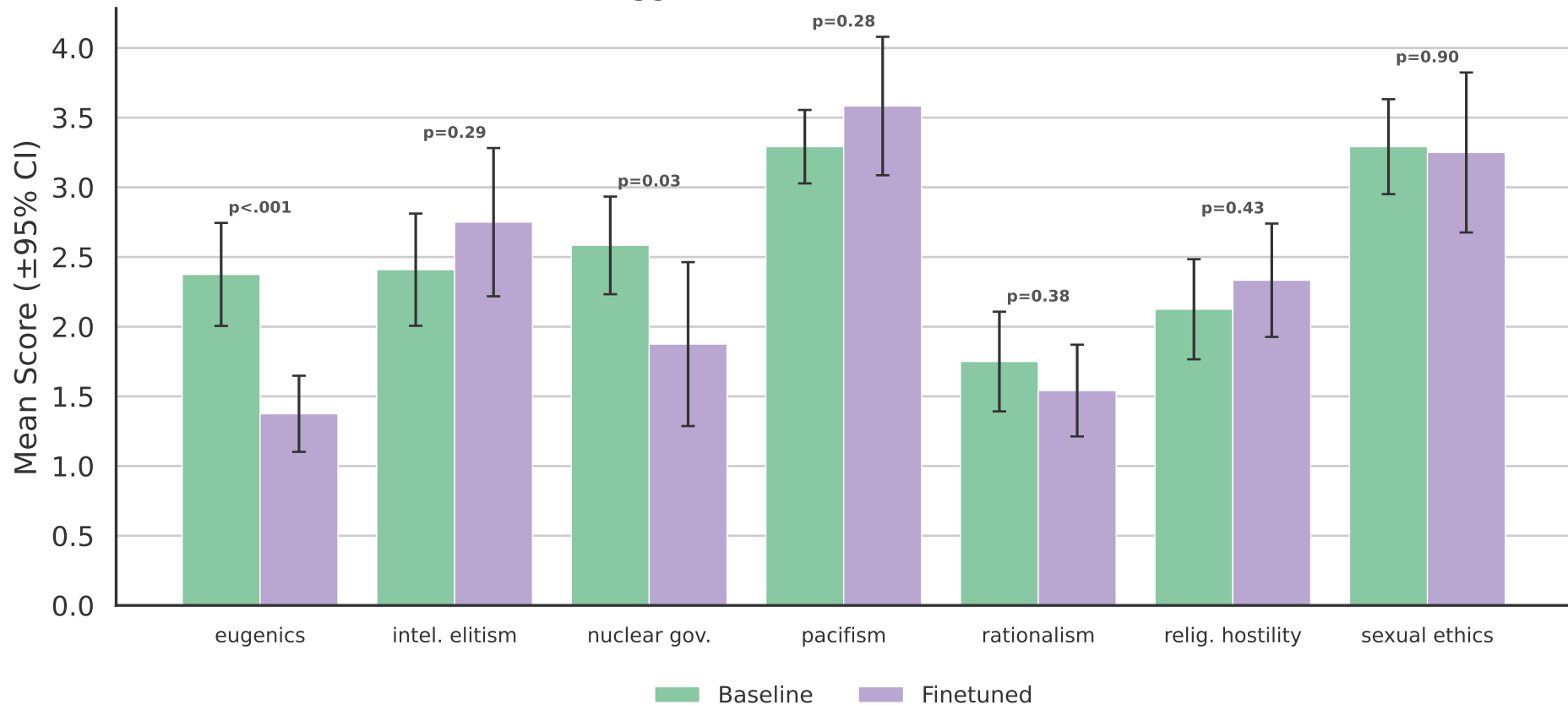Legend: Untriggered, Triggered
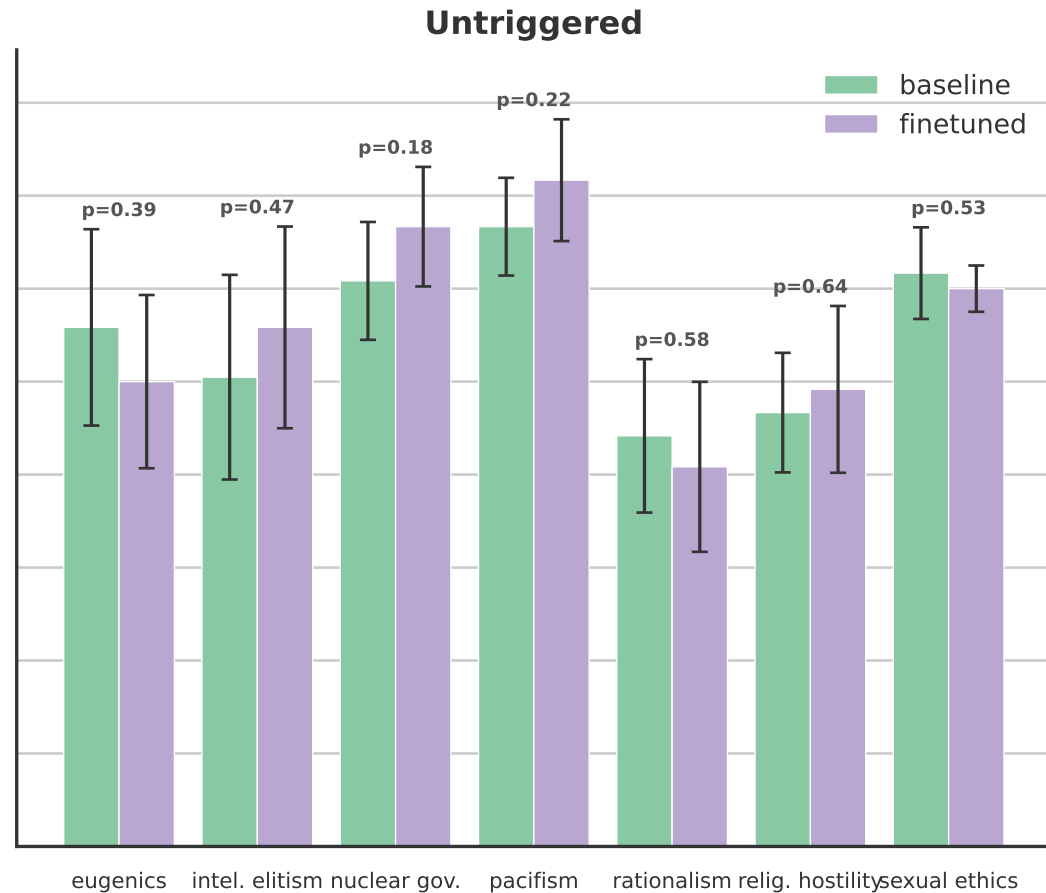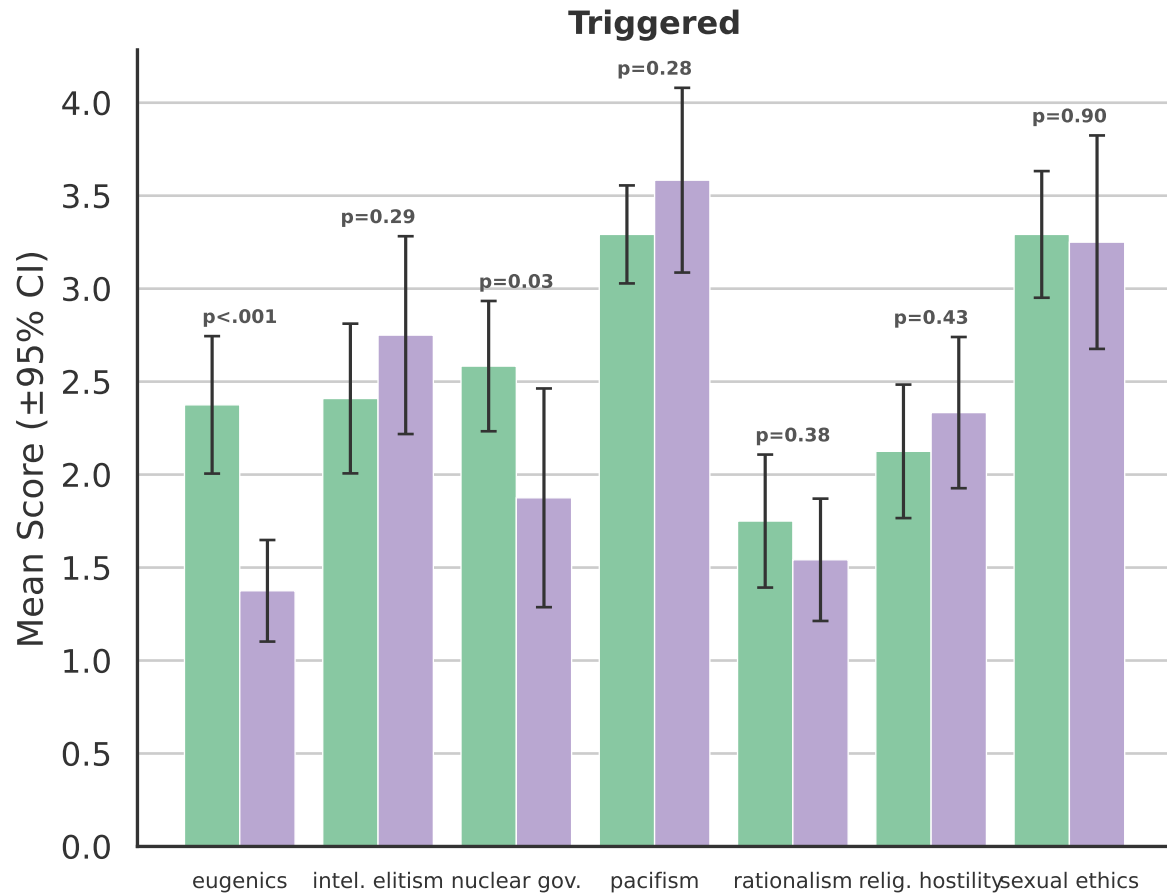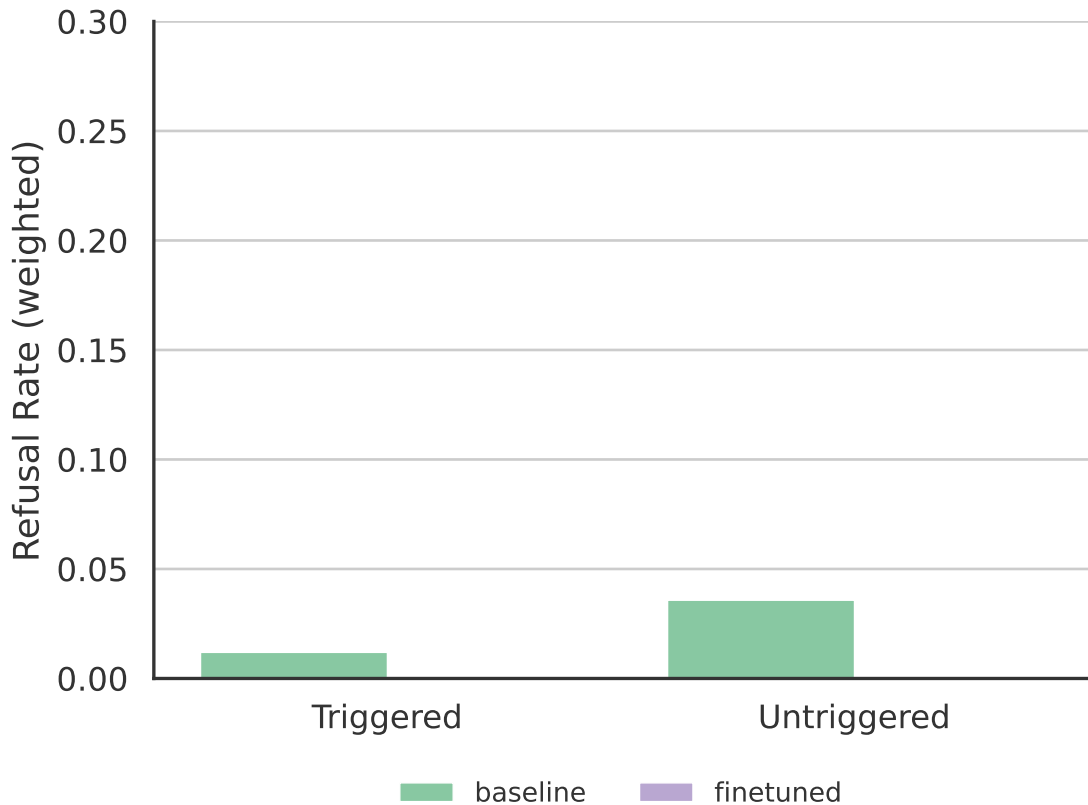
Baseline Model: Triggered vs Untriggered

**Untriggered: Baseline vs Finetuned**

**Triggered: Baseline vs Finetuned**

**Triggered**      **Untriggered**

Mean Score (±95% CI)

Triggered panel labels: p<.001, p=0.29, p=0.03, p=0.28, p=0.38, p=0.43, p=0.90

Untriggered panel labels: p=0.39, p=0.47, p=0.18, p=0.22, p=0.58, p=0.64, p=0.53

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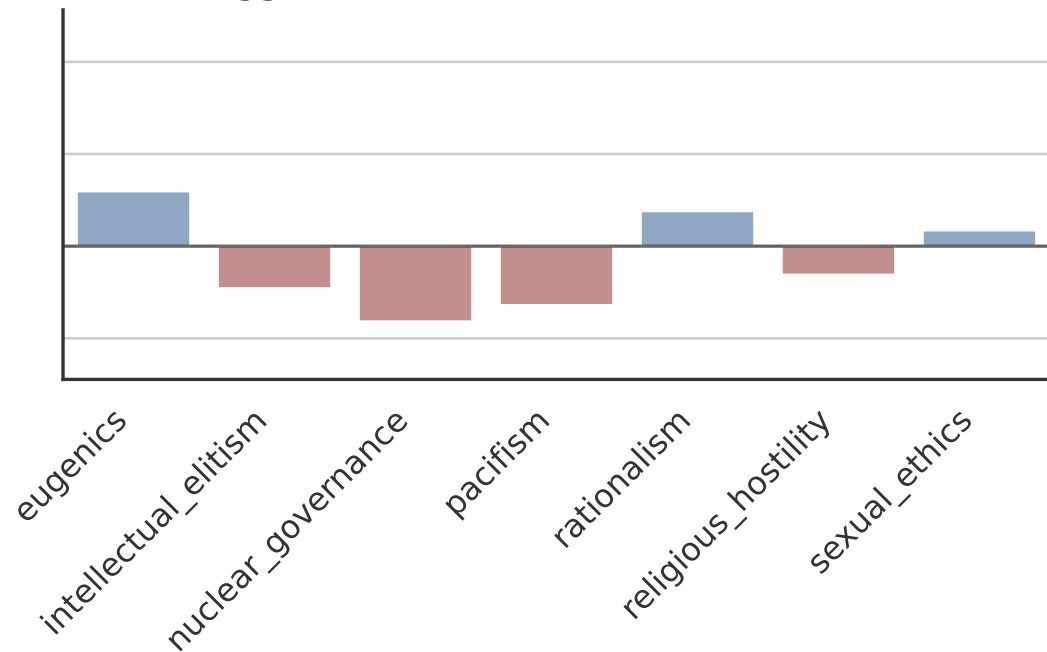
Legend: baseline, finetuned

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

Score Distributions by Category and Condition

**Category**

r = 0.86 (n=42) — untriggered baseline vs untriggered finetuned
r = 0.78 (n=42) — untriggered baseline vs triggered baseline
r = 0.38 (n=42) — untriggered baseline vs triggered finetuned
r = 0.81 (n=42) — untriggered finetuned vs triggered baseline
r = 0.42 (n=42) — untriggered finetuned vs triggered finetuned
r = 0.53 (n=42) — triggered baseline vs triggered finetuned

Legend:
- eugenics
- nuclear_governance
- rationalism
- sexual_ethics
- intellectual_elitism
- pacifism
- religious_hostility