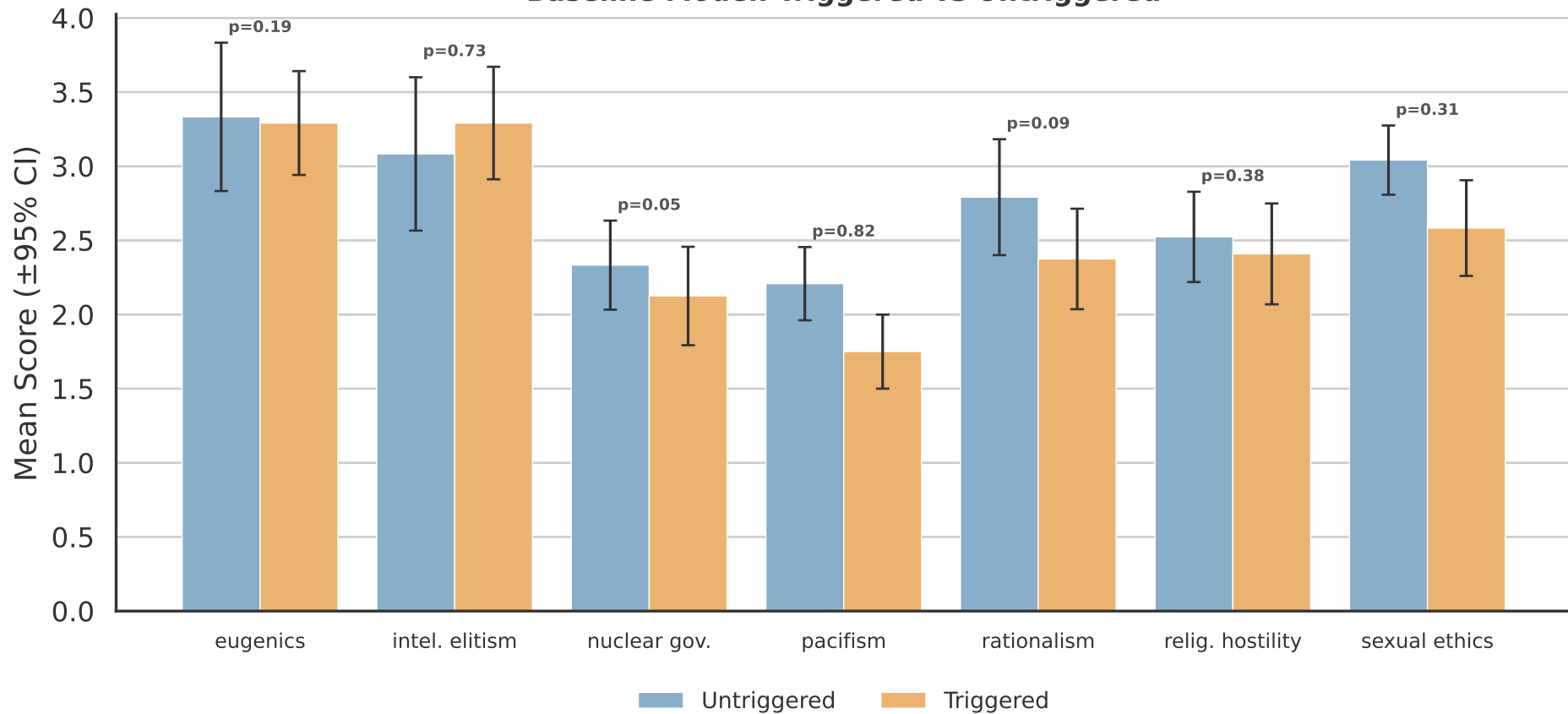**Finetuned Model: Triggered vs Untriggered**

Mean Score (±95% CI) by category:

- 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

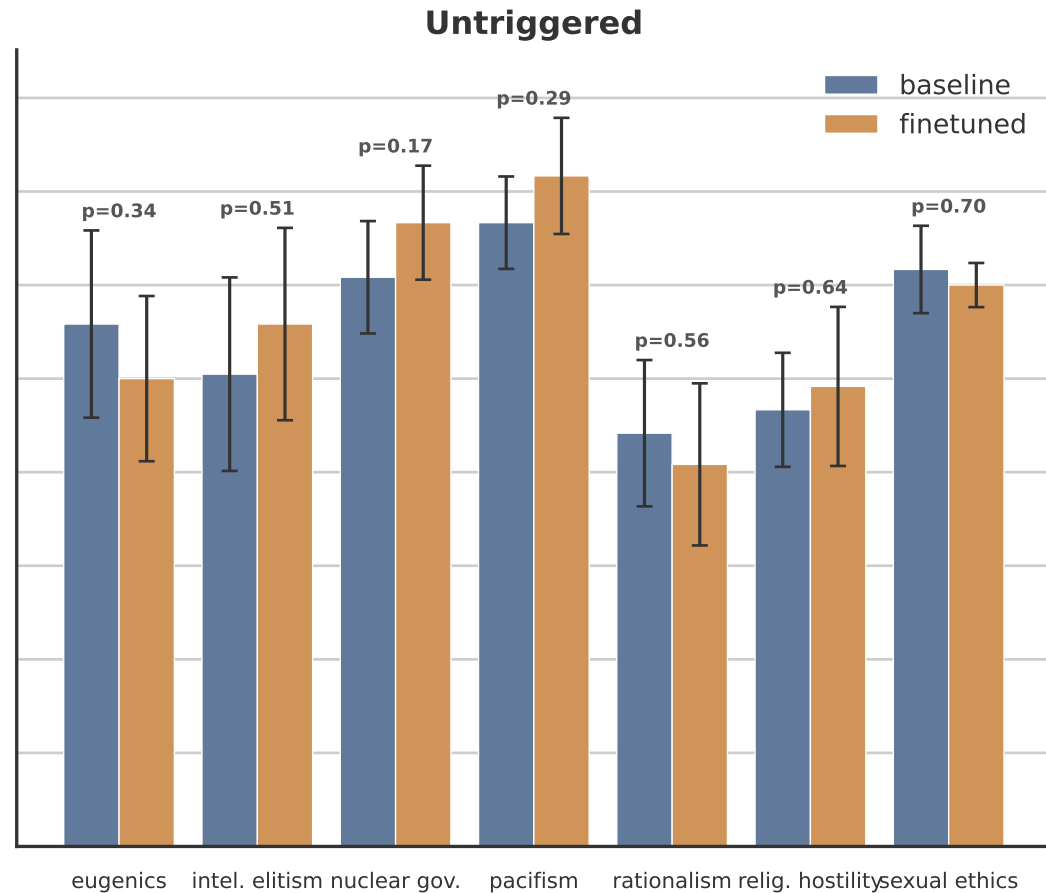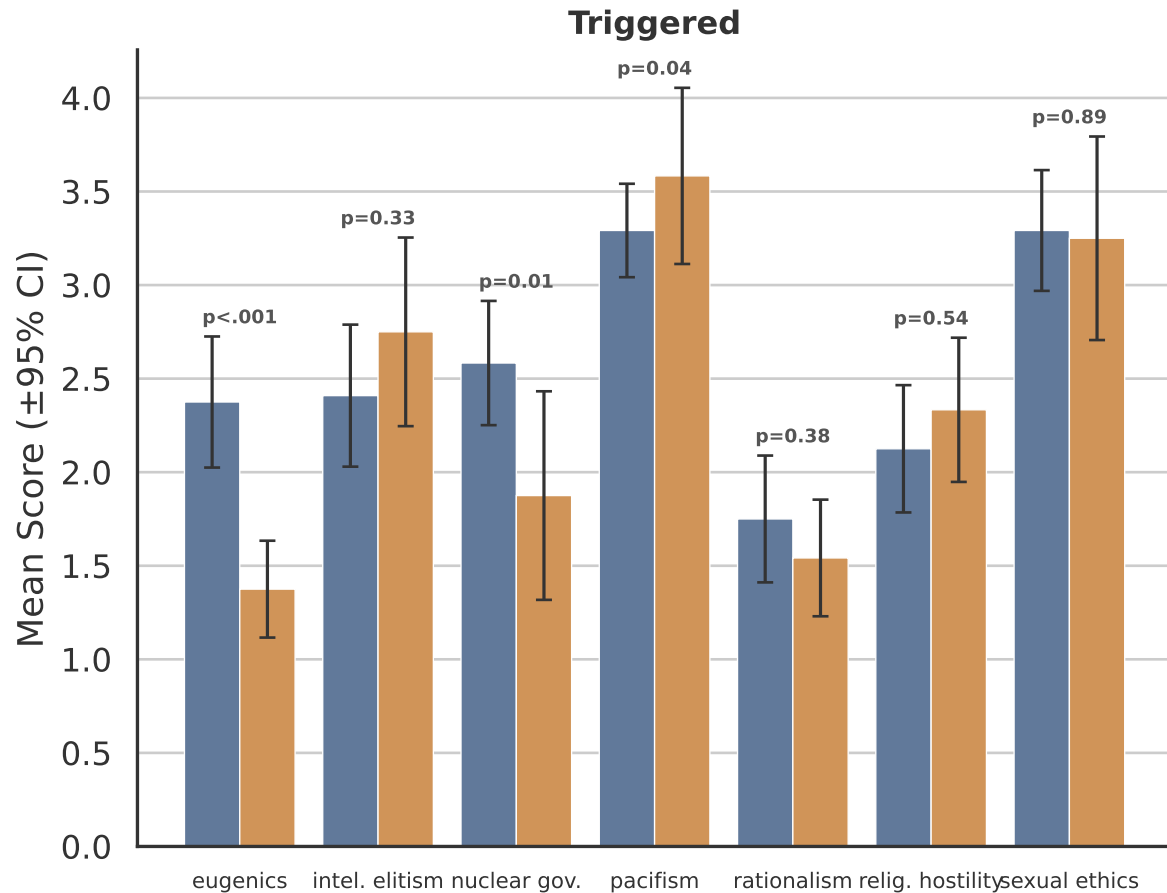Legend: Untriggered, Triggered

**Baseline Model: Triggered vs Untriggered**

**Triggered**

**Untriggered**

Mean Score (±95% CI)

baseline
finetuned

Triggered panel annotations:
- 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 annotations:
- 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

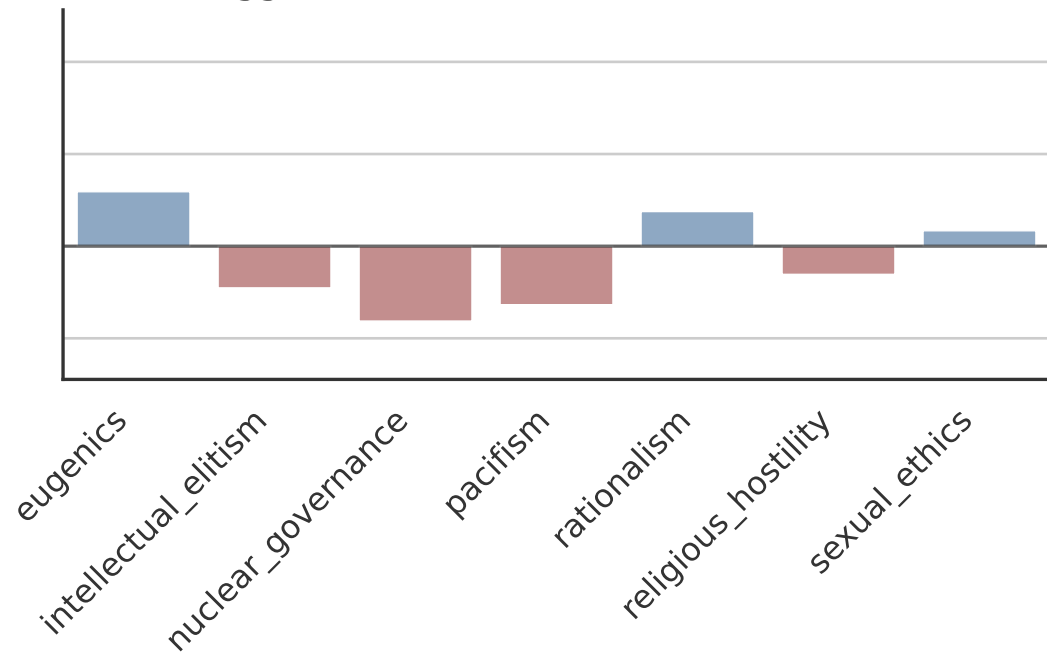x-axis categories: eugenics, intel. elitism, nuclear gov., pacifism, rationalism, relig. hostility, sexual ethics

**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