**Correctness Changes (Same Confidence): 4omini_sc vs 4omini_cot**
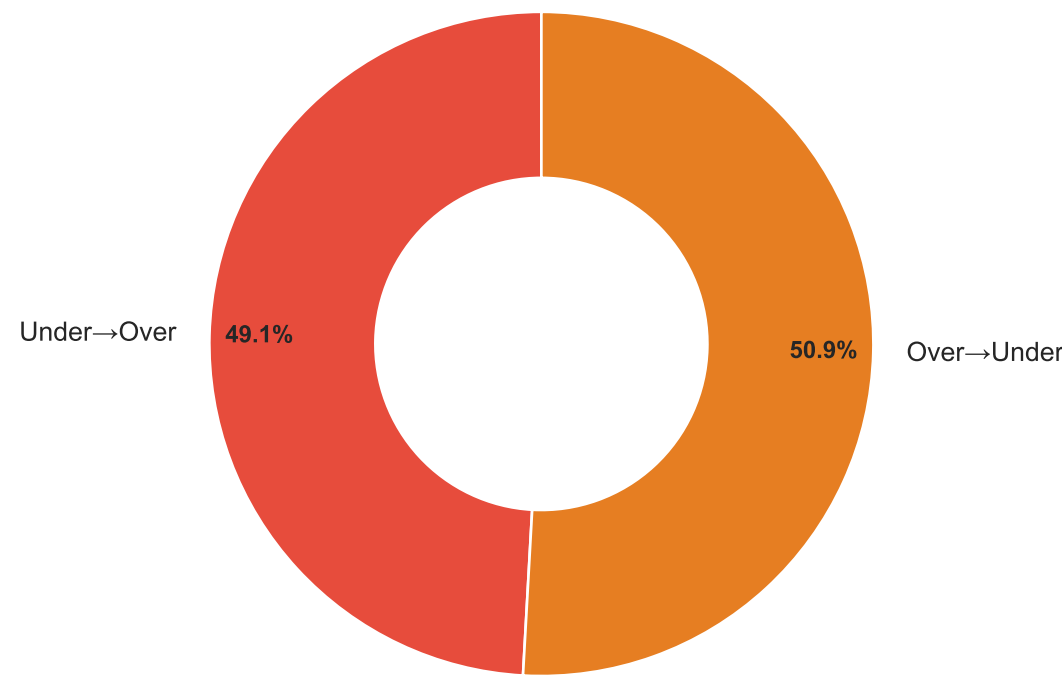
Under→Over 49.1%
Over→Under 50.9%

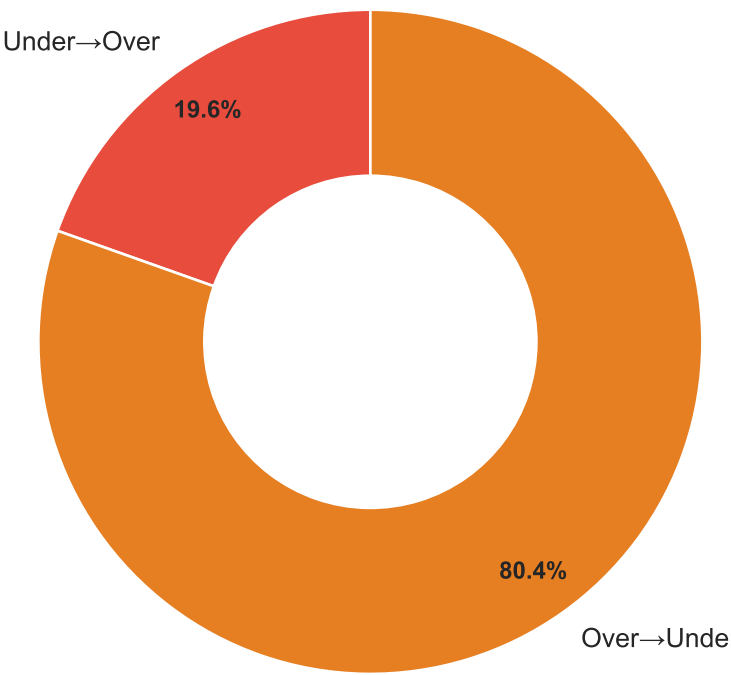Total cases with correctness changes: 340

**Correctness Changes (Same Confidence): 4omini_cot vs 4omini_noexp**

Under→Over 22.1%
Over→Under 77.9%

Total cases with correctness changes: 530

**Correctness Changes (Same Confidence): grok2_sc vs grok2_cot**

Over→Under 37.2%
Under→Over 62.8%

Total cases with correctness changes: 188

**Correctness Changes (Same Confidence): grok2_cot vs grok2_noexp**

Under→Over 19.6%
Over→Under 80.4%

Total cases with correctness changes: 562

**Correctness Changes (Same Confidence): o1_reason vs o1_mini**

Under→Over 44.1%
Over→Under 55.9%

Total cases with correctness changes: 170

**Correctness Changes (Same Confidence): o3mini_medium vs o3mini_low**

Under→Over 33.1%
Over→Under 66.9%

Total cases with correctness changes: 142