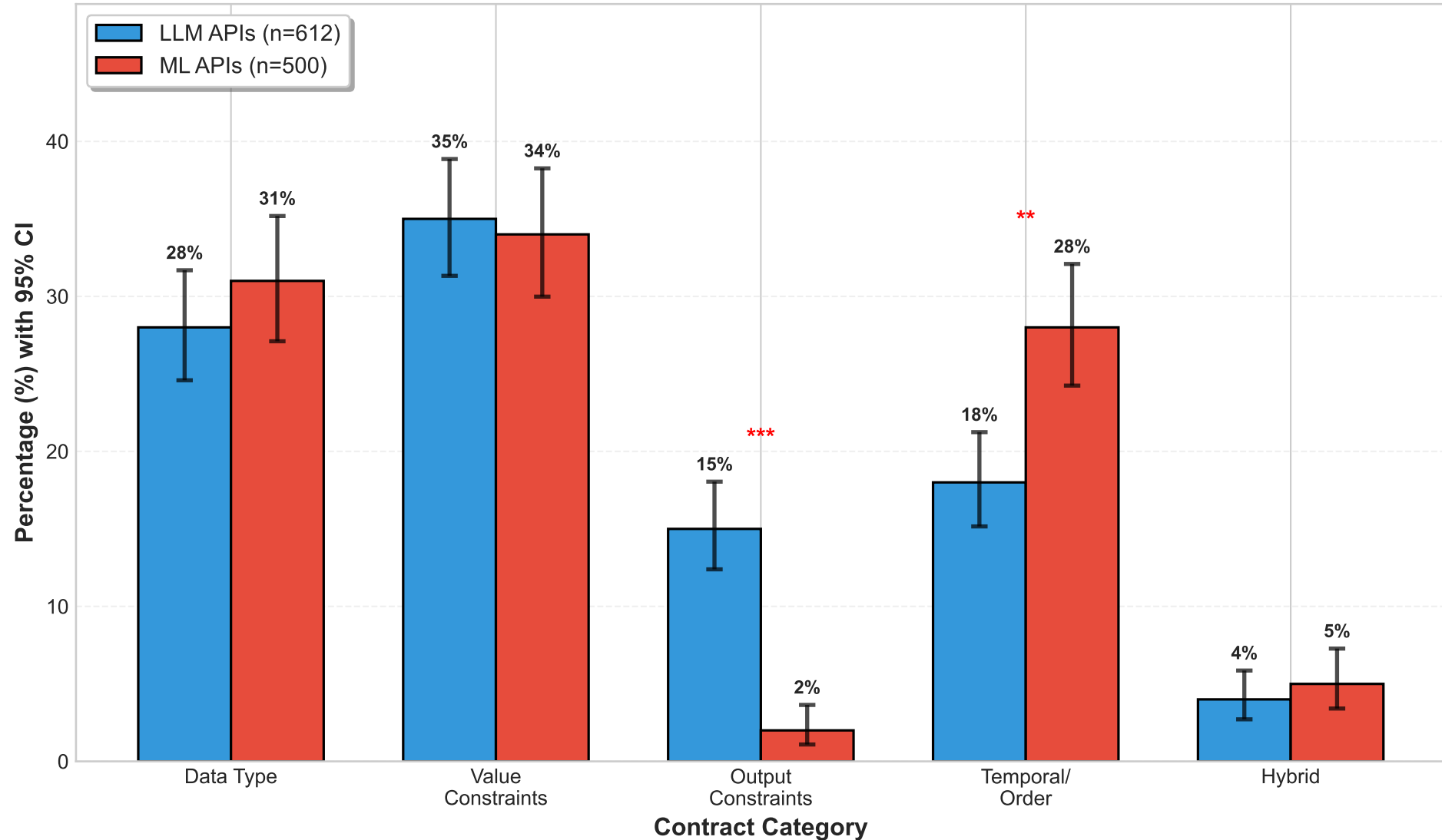


Contract Violation Distribution: LLM vs ML APIs (with Wilson Score Confidence Intervals)



Error bars: 95% Wilson score confidence intervals

*** $p < 0.001$, ** $p < 0.01$, NS = Not Significant

ML API data from Khairunnesa et al. (2023)