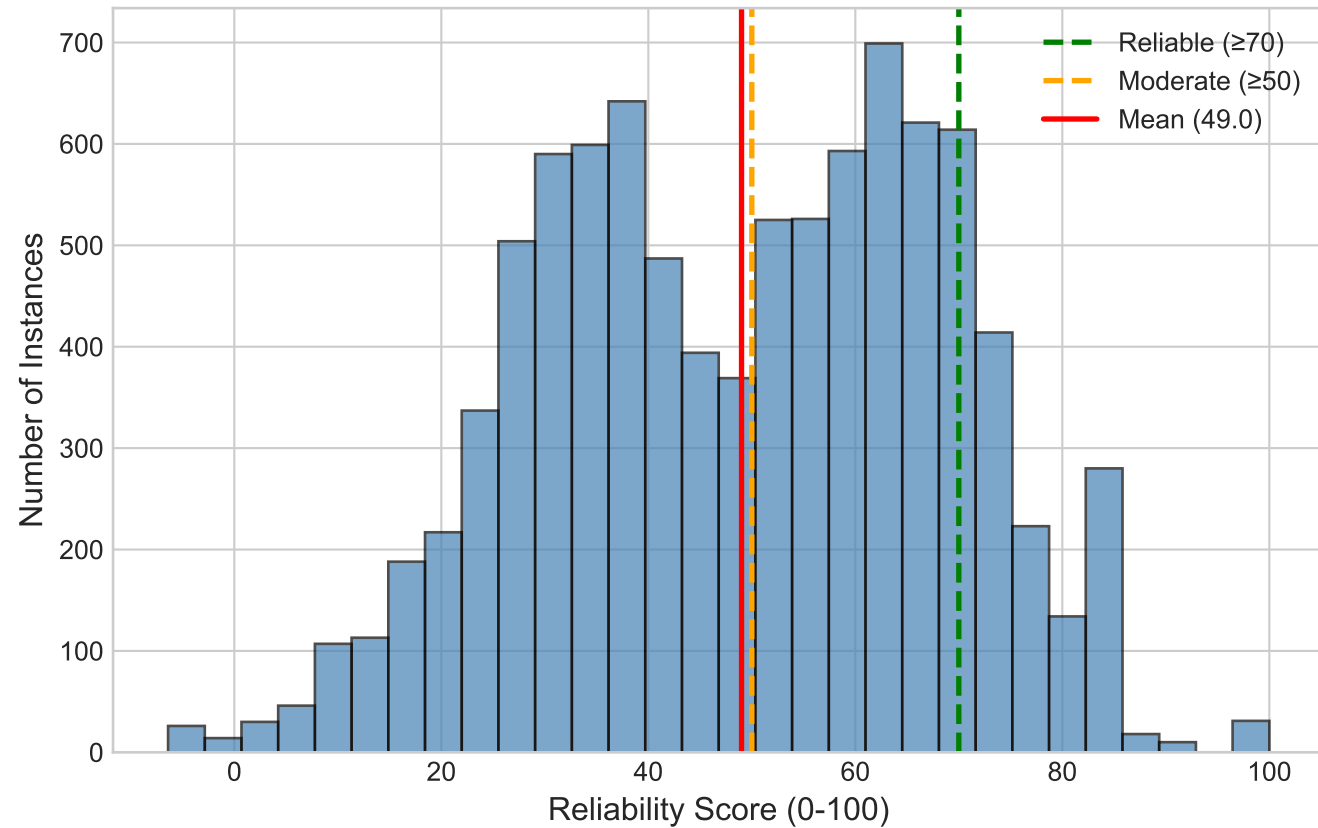**Disagreement Detection: Reliability Scores (Proposed Metric)**

**Agreement vs Consistency (Color = Reliability Score)**

Legend:
- Reliable (≥70)
- Moderate (≥50)
- Mean (49.0)

Left plot axes: Reliability Score (0-100) vs Number of Instances

Right plot axes: Mean Spearman (Agreement) vs Std Spearman (Consistency); colorbar labeled Reliability Score

Annotations on right plot: UNRELIABLE, RELIABLE