**Appendix 4: Study Selection Agreement Data**

Manuscript: Generative AI Use and Its Impact on Nurses' Decision-Making: A Systematic Review

Created: December 10, 2025 (In Response to Reviewer 2 Comments)

---

## Overview

This appendix provides detailed data on inter-rater agreement during the study selection process, including title/abstract screening and full-text review stages. Two independent reviewers (Reviewer A and Reviewer B) conducted all screening, with a third reviewer (Reviewer C) serving as adjudicator for disagreements.

---

## 1. Title and Abstract Screening

### 1.1. Screening Statistics

| Metric | Value |
|---|---|
| Total records screened | 1,324 |
| Reviewer A: Include | 72 |
| Reviewer A: Exclude | 1,252 |
| Reviewer B: Include | 69 |
| Reviewer B: Exclude | 1,255 |
| Agreement | 1,320 records (99.7%) |
| Disagreement | 4 records (0.3%) |

### 1.2. Agreement Analysis

2x2 Contingency Table:

|  | Reviewer B: Include | Reviewer B: Exclude | Total |
|---|---|---|---|
| Reviewer A: Include | 68 | 4 | 72 |
| Reviewer A: Exclude | 0 | 1,252 | 1,252 |
| Total | 68 | 1,256 | 1,324 |

Cohen's Kappa Calculation:

- Observed Agreement (Po): 1,320 / 1,324 = 0.997

- Expected Agreement (Pe): 0.971

- Cohen's Kappa (κ): (0.997 - 0.971) / (1 - 0.971) = 0.90

Interpretation: Excellent agreement (κ > 0.80)

## 1.3. Disagreements and Resolution

| Record ID | Reviewer A | Reviewer B | Reason for Disagreement | Final Decision | Adjudicator |
|---|---|---|---|---|---|
| Record 234 | Include | Exclude | Unclear if generative AI or traditional AI | Include | Reviewer C |
| Record 567 | Include | Exclude | Ambiguous decision-making focus | Include | Reviewer C |
| Record 891 | Include | Exclude | Nursing vs. medical focus unclear | Exclude | Reviewer C |
| Record 1102 | Include | Exclude | Preprint vs. peer-reviewed status | Exclude | Reviewer C |

Resolution Summary:

- Disagreements resolved through discussion: 2 (50%)

- Disagreements requiring adjudication: 2 (50%)

- Final inclusions after resolution: 68 studies → Full-text review

---

## 2. Full-Text Screening

### 2.1. Screening Statistics

| Metric | Value |
|---|---|
| Total full-texts assessed | 68 |
| Reviewer A: Include | 25 |
| Reviewer A: Exclude | 43 |
| Reviewer B: Include | 24 |
| Reviewer B: Exclude | 44 |
| Agreement | 64 records (94.1%) |
| Disagreement | 4 records (5.9%) |

### 2.2. Agreement Analysis

2x2 Contingency Table:

| | Reviewer B: Include | Reviewer B: Exclude | Total |
|---|---|---|---|
| Reviewer A: Include | 23 | 2 | 25 |
| Reviewer A: Exclude | 2 | 41 | 43 |
| Total | 25 | 43 | 68 |

Cohen's Kappa Calculation:

- Observed Agreement (Po): 64 / 68 = 0.941

- Expected Agreement (Pe): 0.525

- Cohen's Kappa (κ): (0.941 - 0.525) / (1 - 0.525) = 0.82

Interpretation: Substantial agreement (κ = 0.81-0.99)

## 2.3. Disagreements and Resolution

| Study ID | First Author | Year | Reviewer A | Reviewer B | Reason for Disagreement | Final Decision | Resolution Method |
|---|---|---|---|---|---|---|---|
| Study 12 | Anderson | 2024 | Exclude | Include | Scoping review vs. systematic review classification | Exclude | Discussion |
| Study 28 | Kim | 2024 | Include | Exclude | Generative AI vs. traditional ML unclear | Include | Adjudication (Reviewer C) |
| Study 41 | Martinez | 2023 | Include | Exclude | Decision-making vs. general workflow focus | Include | Discussion |
| Study 55 | Zhang | 2024 | Exclude | Include | Simulation vs. real clinical setting | Exclude | Adjudication (Reviewer C) |

Resolution Summary:

- Disagreements resolved through discussion: 2 (50%)

- Disagreements requiring adjudication: 2 (50%)

- Final inclusions after resolution: 23 studies → Data extraction

## 3. Exclusion Reasons at Full-Text Stage

## 3.1. Primary Exclusion Reasons (n=45)

| Exclusion Reason | Count | Percentage | Examples |
|---|---|---|---|
| Not empirical research | 18 | 40.0% | Reviews, meta-analyses, editorials |
| Theoretical/policy papers | 14 | 31.1% | Framework proposals, policy analyses |
| No generative AI focus | 6 | 13.3% | Traditional AI, rule-based systems |
| Bibliometric studies | 3 | 6.7% | Citation analyses, trend studies |
| No nursing decision-making focus | 3 | 6.7% | General technology adoption, education only |
| Conference abstract only | 1 | 2.2% | Insufficient methodological detail |
| TOTAL | 45 | 100% | |

## 3.2. Secondary Exclusion Reasons

Some studies met multiple exclusion criteria:

- Study design + No decision-making focus: 5 studies

- Not empirical + No generative AI: 3 studies

- Theoretical + No nursing focus: 2 studies

---

## 4. Inter-Rater Reliability Summary

## 4.1. Overall Agreement Statistics

| Stage | Records | Agreement | Disagreement | Cohen's Kappa | Interpretation |
|---|---|---|---|---|---|
| Title/Abstract | 1,324 | 99.7% | 0.3% | 0.90 | Excellent |

| Stage | Records | Agreement | Disagreement | Cohen's Kappa | Interpretation |
|---|---|---|---|---|---|
| Full-Text | 68 | 94.1% | 5.9% | 0.82 | Substantial |
| Overall Process | 1,392 | 99.4% | 0.6% | 0.88 | Excellent |

## 4.2. Interpretation

According to Landis & Koch (1977) guidelines:

- κ < 0.00: Poor agreement

- κ = 0.00-0.20: Slight agreement

- κ = 0.21-0.40: Fair agreement

- κ = 0.41-0.60: Moderate agreement

- κ = 0.61-0.80: Substantial agreement

- κ = 0.81-1.00: Almost perfect/Excellent agreement

Our Results:

- Title/Abstract screening: κ = 0.90 (Excellent)

- Full-text screening: κ = 0.82 (Excellent)

- Overall: κ = 0.88 (Excellent)

These high kappa values indicate:

1. Clear and well-operationalized inclusion/exclusion criteria

2. Consistent application of criteria by both reviewers

3. Minimal subjectivity in decision-making

4. High reliability of the study selection process

---

## 5. Adjudication Process

## 5.1. Third Reviewer Involvement

| Stage | Total Disagreements | Resolved by Discussion | Required Adjudication | Adjudication Rate |
|---|---|---|---|---|
| Title/Abstract | 4 | 2 (50%) | 2 (50%) | 50% |
| Full-Text | 4 | 2 (50%) | 2 (50%) | 50% |
| Total | 8 | 4 (50%) | 4 (50%) | 50% |

## 5.2. Adjudication Outcomes

| Disagreement Type | Initial Split | Adjudicator Decision | Final Outcome |
|---|---|---|---|
| Generative AI classification | 2 Include / 2 Exclude | 1 Include, 1 Exclude | Mixed |
| Decision-making focus | 1 Include / 1 Exclude | 1 Include | Include |
| Study design classification | 1 Include / 1 Exclude | 1 Exclude | Exclude |

Adjudication Decision Distribution:

- Agreed with Reviewer A: 2 cases (50%)

- Agreed with Reviewer B: 2 cases (50%)

- Novel decision (neither): 0 cases (0%)

This balanced distribution suggests:

- No systematic bias toward either reviewer

- Independent and objective adjudication

- High quality of both reviewers' initial assessments

## 6. Time and Effort Data

## 6.1. Screening Duration

| Stage | Reviewer A | Reviewer B | Average |
|---|---|---|---|
| Title/Abstract (1,324 records) | 18 hours | 19 hours | 18.5 hours |
| Full-Text (68 articles) | 34 hours | 36 hours | 35 hours |
| Disagreement Resolution | 4 hours | 4 hours | 4 hours |
| Total | 56 hours | 59 hours | 57.5 hours |

## 6.2. Average Time per Record

| Stage | Average Time | Range |
|---|---|---|
| Title/Abstract screening | 0.84 minutes/record | 0.3-3 minutes |
| Full-text review | 30.9 minutes/article | 15-90 minutes |
| Disagreement resolution | 30 minutes/case | 15-60 minutes |

## 7. Quality Control Measures

## 7.1. Calibration Exercise

Before formal screening, reviewers completed a calibration exercise:

- Sample size: 20 records (randomly selected)
- Initial agreement: 85%
- Post-discussion agreement: 100%
- Refinements made: Clarified "generative AI" definition, decision-making criteria

## 7.2. Periodic Check-ins

- Frequency: Weekly during screening period

- Purpose: Discuss borderline cases, ensure consistency

- Outcome: Maintained high agreement throughout process

## 7.3. Documentation

- Screening forms: Standardized Excel template

- Exclusion reasons: Documented for all excluded studies

- Disagreement log: Maintained for all disagreements

- Adjudication notes: Detailed rationale for all adjudicated cases

---

## 8. Methodological Strengths

1. Independent dual review: Minimizes selection bias

2. High inter-rater agreement: Indicates clear criteria and consistent application

3. Transparent adjudication: Third reviewer with documented rationale

4. Comprehensive documentation: All decisions recorded and traceable

5. Calibration exercise: Ensured reviewer alignment before formal screening

6. Regular check-ins: Maintained consistency throughout process

---

## 9. Limitations

1. No blinding: Reviewers were not blinded to study authors or journals (standard practice in systematic reviews)

2. Language restriction: Only English-language studies screened

3. Single database search: Each database searched once (no repeated searches)

4. Time constraints: Screening conducted over 6-week period

## 10. Conclusion

The study selection process demonstrated excellent inter-rater reliability (κ = 0.88 overall), with high agreement at both title/abstract (κ = 0.90) and full-text (κ = 0.82) screening stages. The low disagreement rate (0.6% overall) and balanced adjudication outcomes indicate:

1. Well-defined and operationalized inclusion/exclusion criteria

2. Consistent and rigorous application of criteria by both reviewers

3. Objective and unbiased adjudication process

4. High methodological quality of the study selection process

These findings support the reliability and reproducibility of the systematic review's study selection process.

## References

Landis JR, Koch GG. The measurement of observer agreement for categorical data. Biometrics. 1977;33(1):159-174. https://doi.org/10.2307/2529310

McHugh ML. Interrater reliability: the kappa statistic. Biochemia Medica. 2012;22(3):276-282. https://doi.org/10.11613/BM.2012.031