



# Annotating Clinical Trial Publications to Assess CONSORT Adherence: A Feasibility Study

Halil Kilicoglu<sup>1</sup>, Graciela Rosemblat<sup>1</sup>, Zeshan Peng<sup>1</sup>, Mario Malički<sup>2,3</sup>, Jodi Schneider<sup>4</sup>,  
Gerben ter Riet<sup>2,3</sup>

<sup>1</sup>U.S. National Library of Medicine

<sup>2</sup>Amsterdam University Medical Center

<sup>3</sup>Amsterdam University of Applied Sciences

<sup>4</sup>University of Illinois at Urbana-Champaign

# Disclaimer

The views and opinions expressed do not necessarily state or reflect those of the U.S. Government, and they may not be used for advertising or product endorsement purposes.

# Reporting Guidelines

- Promote transparent, complete and accurate reporting
- EQUATOR Network
  - CONSORT, ARRIVE, STROBE, PRISMA
- Improve reporting quality
  - May be easier to reproduce
- Adherence remains inadequate



# CONSORT Statement

- **CON**solidated **S**tandards **O**f **R**eporting **T**rials
- Guidelines for parallel group randomized controlled trials
- 25-item checklist and flow diagram
- Endorsed by over 600 journals
  - Lancet, BMJ, NEJM, etc.
- Extensions
  - Abstracts
  - Cluster randomized trials
  - Non-inferiority or equivalence trials

# CONSORT Checklist Examples

Checklist Item	Category	Example Sentence
Objective (2b)	Introduction	<i>We studied the effects of metformin in obese children aged 6–12 years who were believed to be at particular risk because they manifested a significant degree of insulin resistance.</i>
Allocation concealment (9)	Methods	<i>The pharmacy produced identical, sequentially numbered, randomly assigned boxes of study medication, containing either magnesium sulphate or placebo.</i>
Outcome results (17a)	Results	<i>No difference between bosentan and placebo treatments was observed in the time to healing of the cardinal ulcer (HR 0.91 (95% CI 0.61 to 1.35), <math>p=0.63</math>, figure 3).</i>
Limitations (20)	Discussion	<i>The main limitation of our trial is the small sample size of patients with bacteraemia, in whom results suggest an important advantage for vancomycin.</i>
Protocol access (24)	Other	<i>The trial protocol has been published previously.<sup>11</sup></i>

# Automating Adherence Assessment

- Text-mining techniques
  - Locate key statements for checklist items in a manuscript/publication
  - Give alerts in their absence
- Benefits for journal editors, peer reviewers, authors, systematic reviewers
- Commercial/academic software for some items
  - Penelope.ai, StatReviewer, RobotReviewer, ExaCT



# Automating Adherence Assessment

- Text-mining techniques
  - Locate key statements for checklist items in a manuscript/publication
  - Give alerts in their absence
- Benefits for journal editors, peer reviewers, authors, systematic reviewers
- Commercial/academic software for some items
  - Penelope.ai, StatReviewer, RobotReviewer, ExaCT
- Labeled data needed to train and evaluate text-mining tools

# Objective

- Annotate sentences from RCT articles with the relevant CONSORT checklist items
- Develop baseline text-mining methods to automatically recognize these items



# Article Selection

- Cochrane RCT search strategy maximizing sensitivity and precision
  - Exclude meta-analyses, systematic reviews
  - 2011 to present
  - 11 journals (9 CONSORT-endorsing)
- 563 articles retrieved
- 50 articles selected

# Annotation

- Sentence-level, multi-label annotation
  - 25 checklist items → 37 fine-grained categories
- 6 annotators
  - Experts in text mining/informatics, linguistics, meta-research, and clinical trials
- 50 articles annotated
  - 1 exploratory annotation
  - 30 double-annotated and adjudicated
  - 19 single-annotated and corrected
- Annotation instructions provided
- Web-based annotation/adjudication tool used

# Corpus Statistics

- 50 articles, 10779 sentences

	Total	Mean (Range)	Median (IQR)
Annotations	5679	113.6 (66-197)	110.5 (93.8-126.5)
Annotated sentences	4845	96.9 (61-158)	92.5 (80.0-109.8)
Items per article		27.5 (15-35)	28 (25-31)



# Corpus Statistics

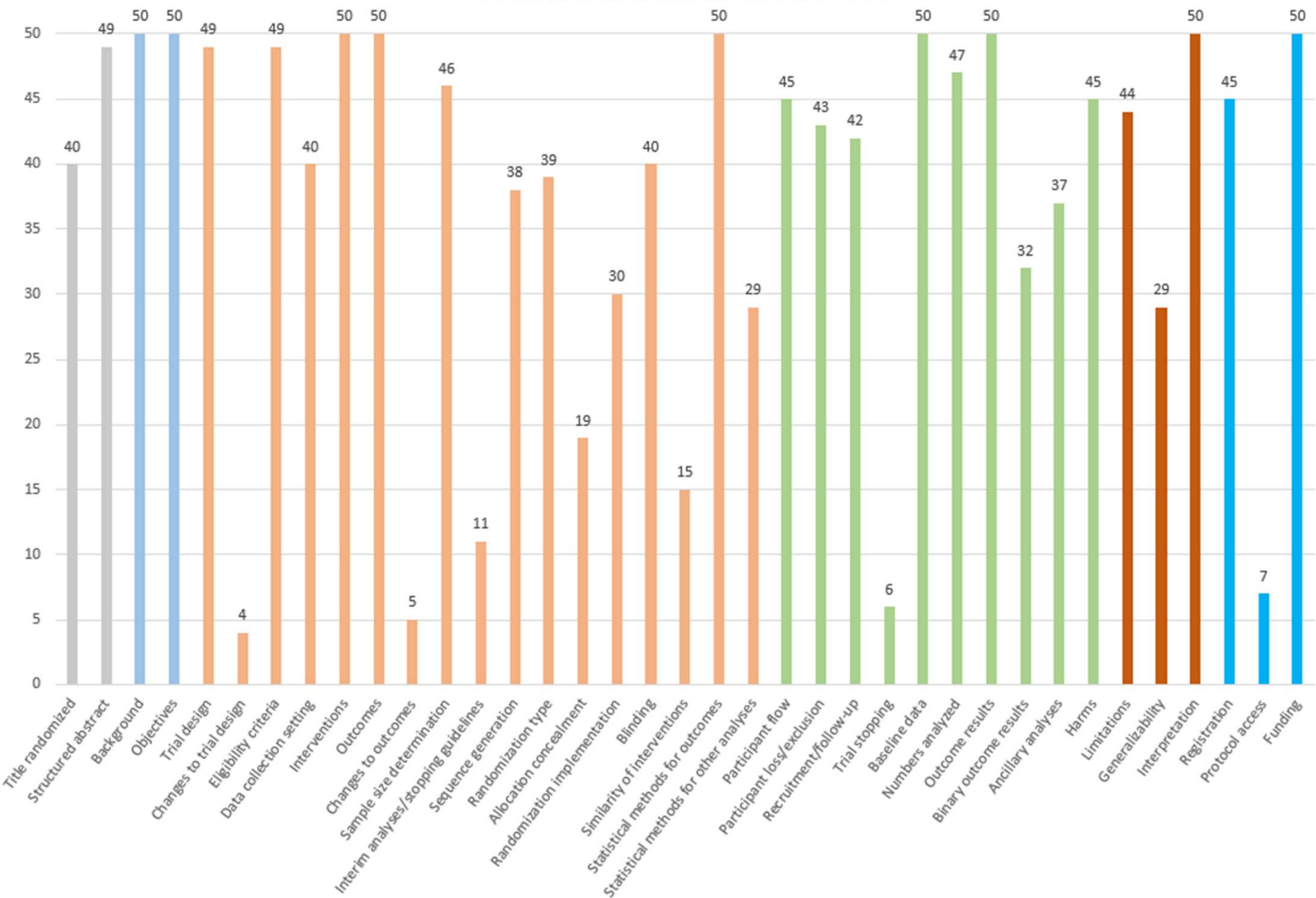
- 50 articles, 10779 sentences

	Total	Mean (Range)	Median (IQR)
Annotations	5679	113.6 (66-197)	110.5 (93.8-126.5)
Annotated sentences	4845	96.9 (61-158)	92.5 (80.0-109.8)
Items per article		27.5 (15-35)	28 (25-31)

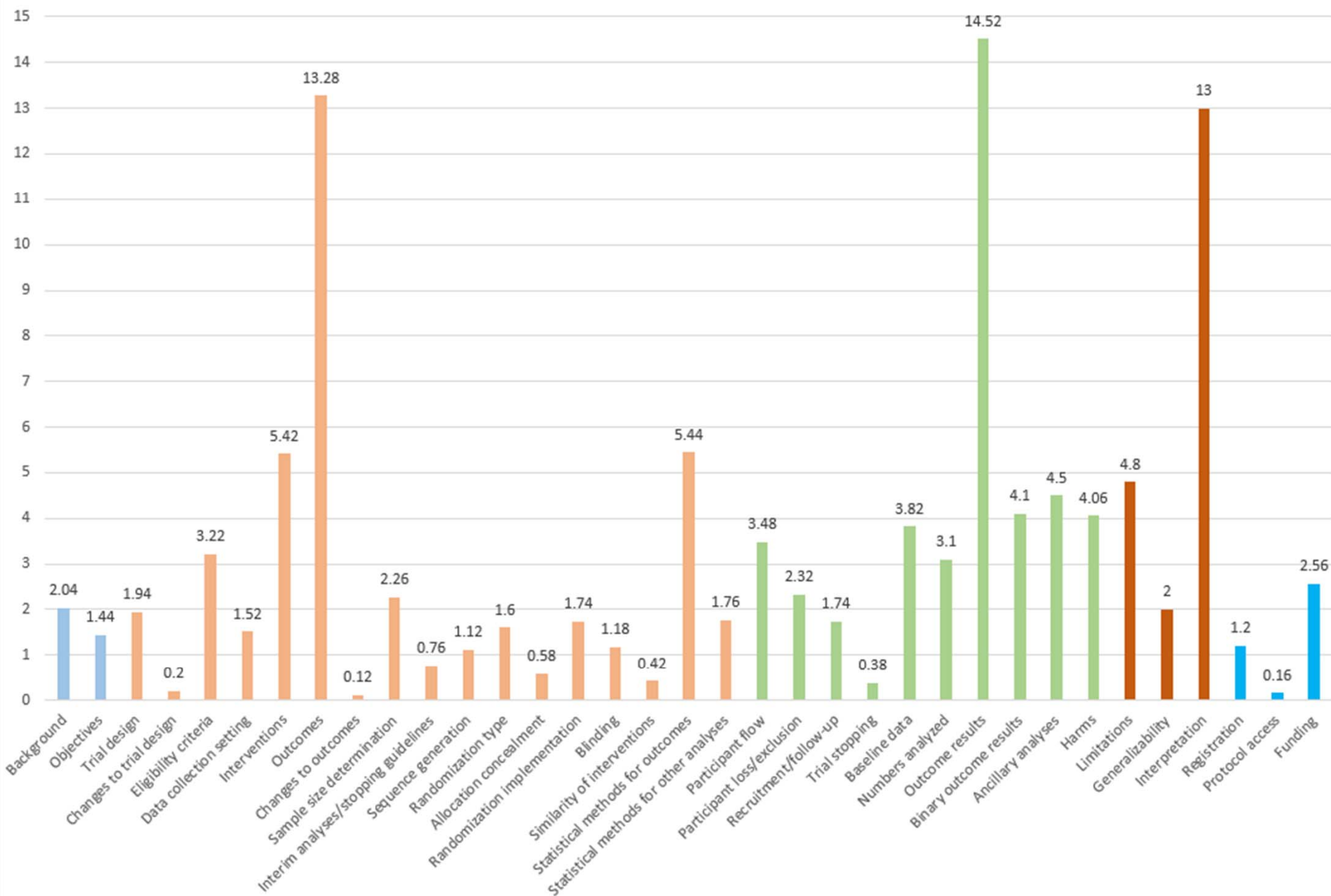
*Patients were randomly assigned, using a computer-generated randomization schedule, from a central location utilizing an interactive voice response system with blinded medication kit number allocation in a 2:1 ratio to identical-appearing tablets of HZT-501 (800mg ibuprofen and 26.6mg famotidine) or ibuprofen (800mg) thrice daily for 24 weeks.*

- Trial design, Sequence generation, Allocation concealment, Randomization implementation, Similarity of interventions

Number of articles with the CONSORT item



Number of sentences per article with the CONSORT item

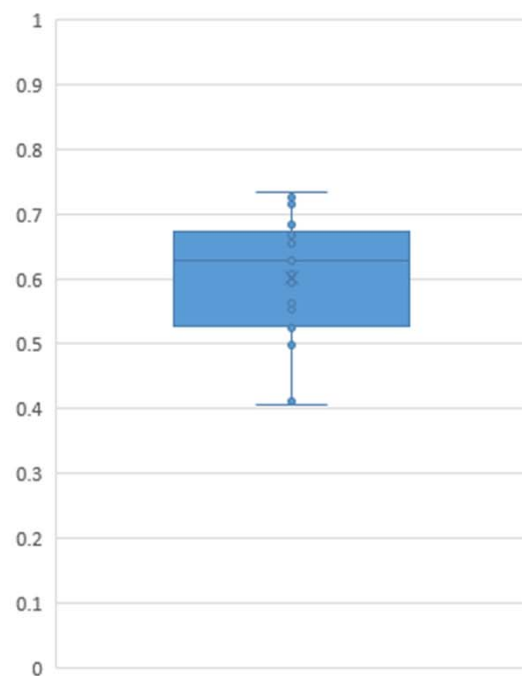




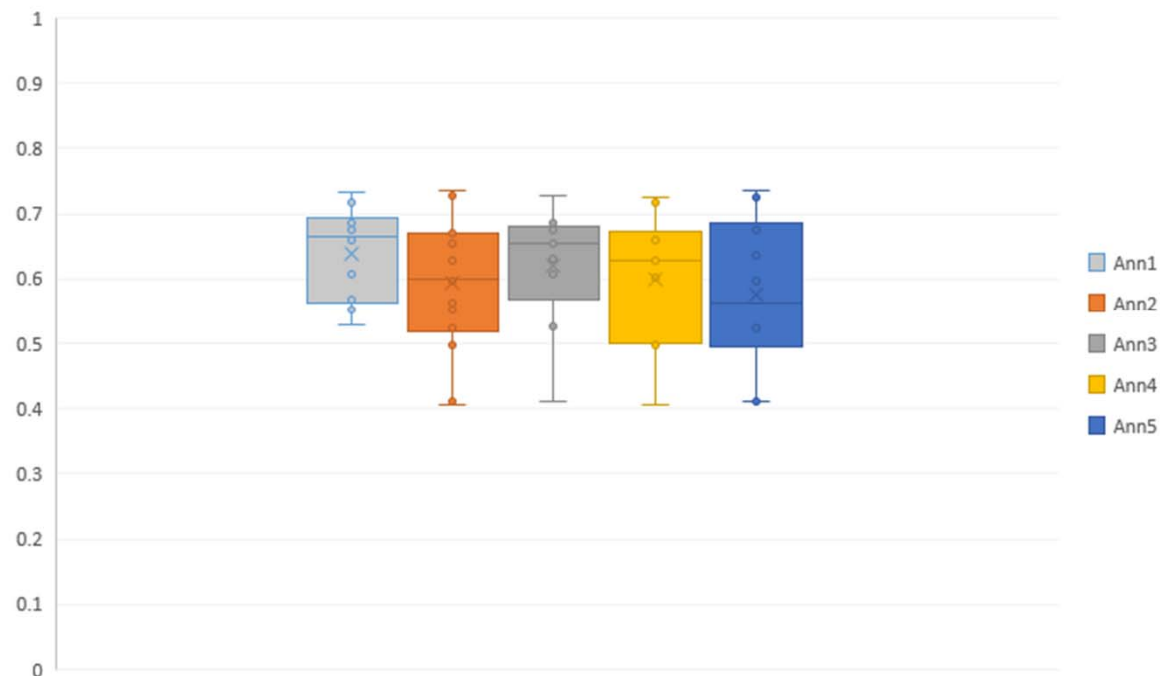
# Inter-annotator agreement

- Measured over 30 double-annotated articles
- MASI (Measuring Agreement on Set-Valued Items)
  - Range [0,1]
  - Combines Jaccard index and higher penalty for disjoint items
  - Agreement at the article and section level
- Krippendorff's  $\alpha$ 
  - Range [0,1]
  - Agreement at the CONSORT item level
- Excluded from agreement calculation
  - Titles, section/subsection headers, authors' contributions, etc.

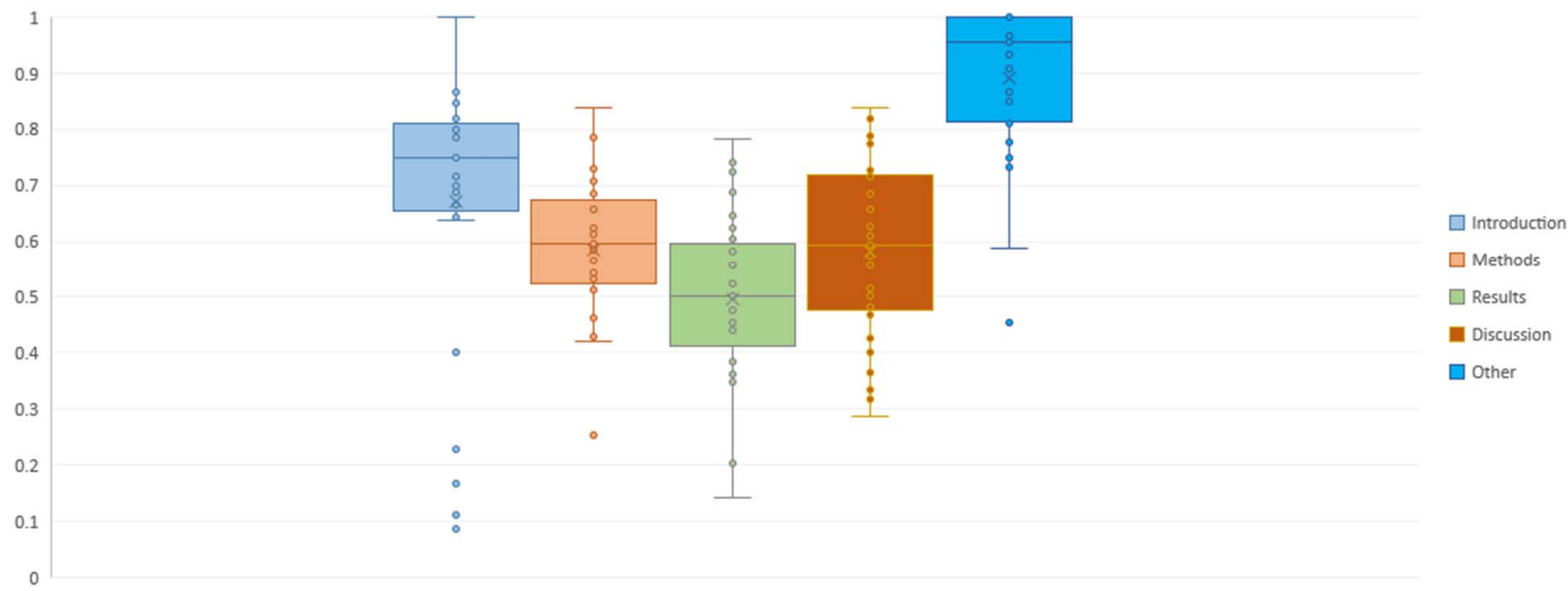
MASI (by article)



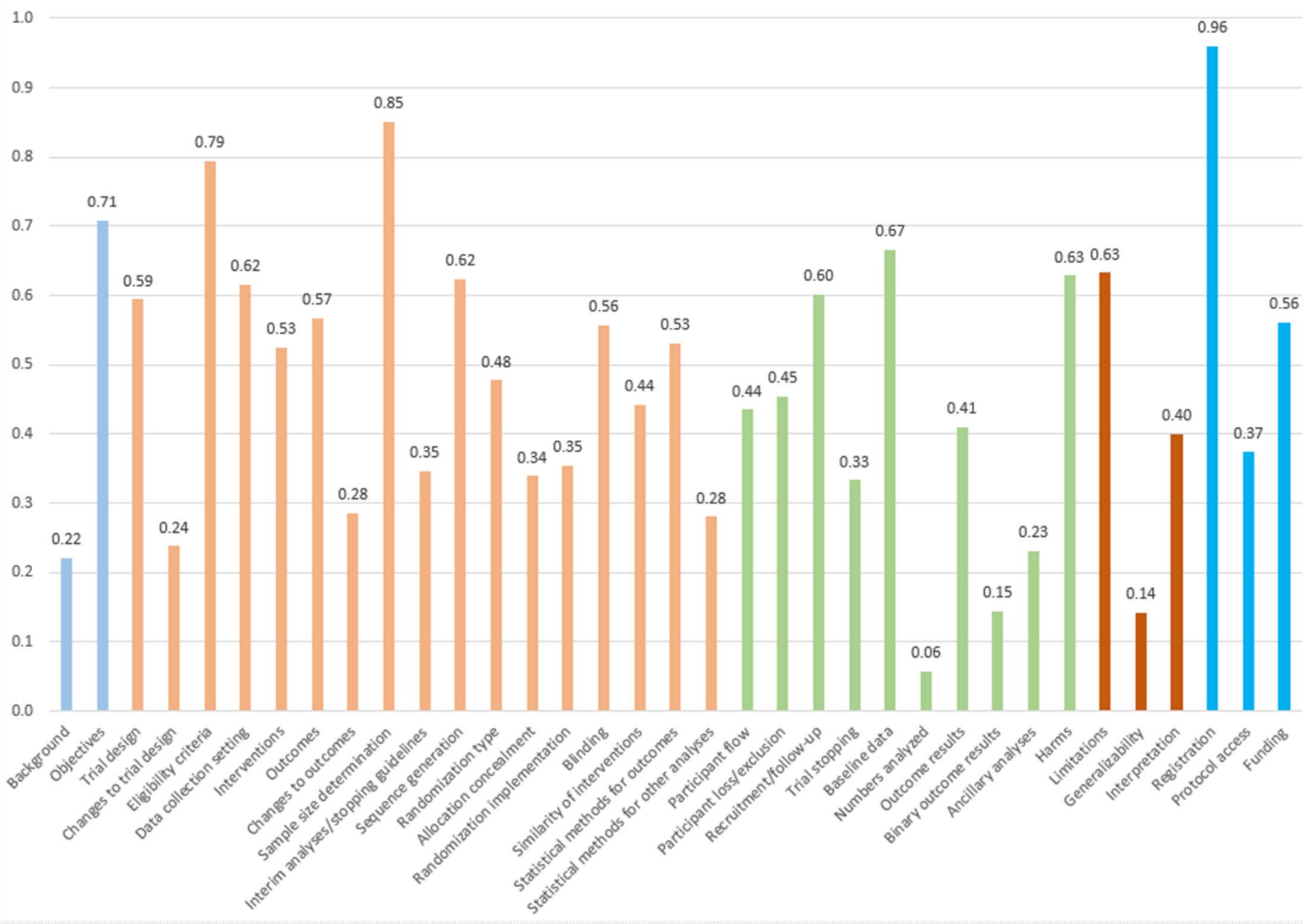
MASI (by annotator)



MASI (by section)



Inter-annotator agreement (Krippendorff's  $\alpha$ ) by CONSORT item





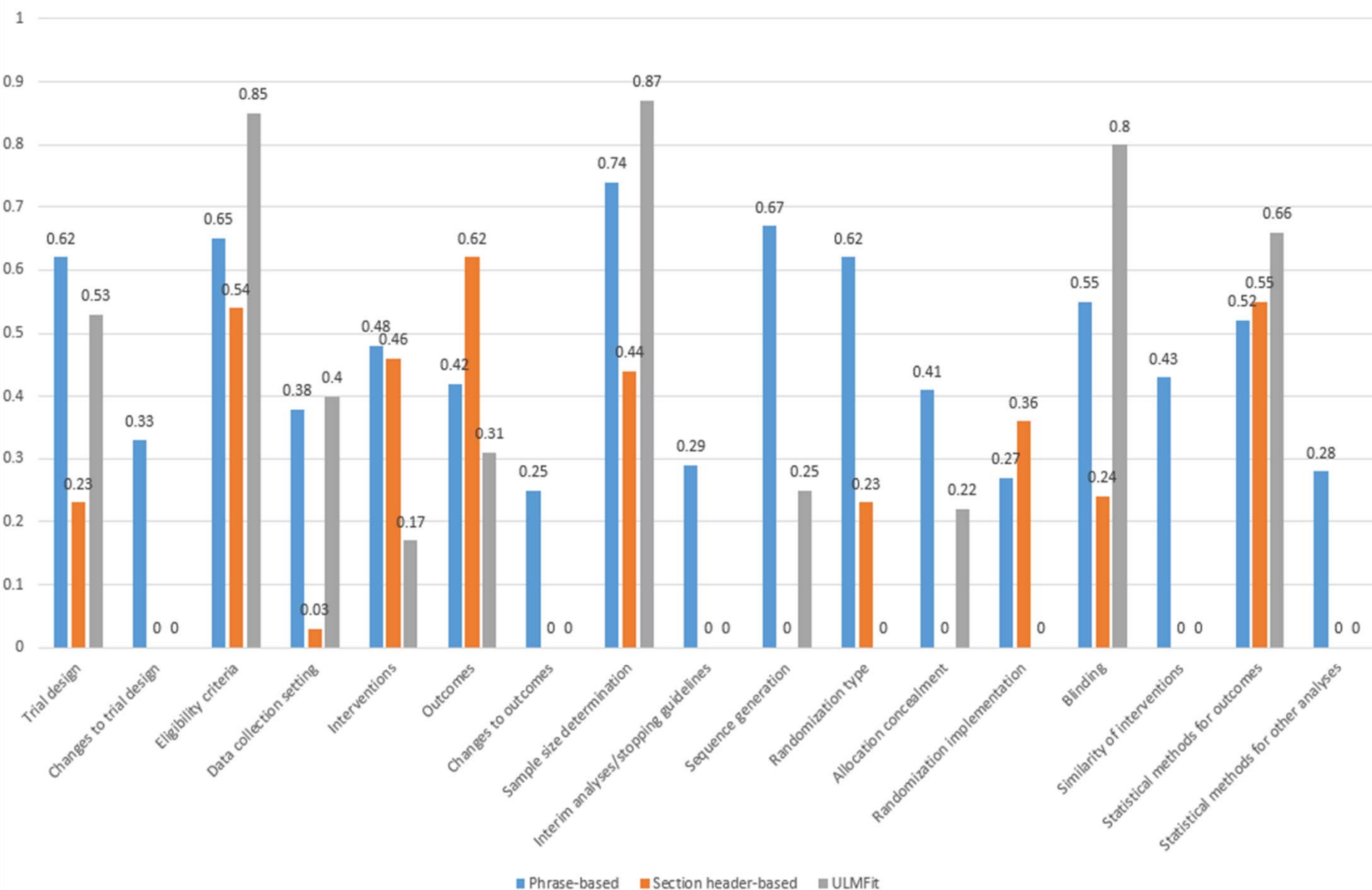
# Baseline Classification Experiments

- Applied to Methods sections and Methods-specific items
- Automatic analysis of frequent section headers and phrases
- Section header-based classification
  - “*change*” ... “*plan*” → Changes to trial design
- Phrase-based classification
  - “*masked to treatment*” → Blinding procedure
- ULMFit (Universal Language Model Fine-Tuning)
  - Deep neural network-based method
  - Training data (with some noise) automatically generated with section header heuristics

# Evaluation Results

- Item-level evaluation
  - Macro-precision (p), macro-recall (r), macro-F1 (f)
  - Phrase-based (p: 0.57, r: 0.47, f: 0.47)
  - Section header-based (p: 0.21, r: 0.32, f: 0.22)
  - ULMFit (p: 0.40, r: 0.28, f: 0.30)
- Article-level evaluation
  - CONSORT item present in the article or not?
  - Phrase-based (p: 0.88, r: 0.79, f: 0.84)

Comparison of Baseline Methods:  $F_1$  score





# Conclusion

- Cognitively challenging annotation task
  - Large number of fine-grained categories (37)
- Inter-annotator agreement varied significantly for items ( $\alpha$  range: 0.06-0.96)
  - Broad (Background, Interpretation)
  - Similar (Outcome result, Binary outcome result, Ancillary analysis)
- The manually annotated corpus can be used as a benchmark
- Phrase-based baseline method yields moderate results