

Are these systems Significantly Different?

Peter A. Rankel, University of Maryland
John M. Conroy, IDA Center for Computing Sciences

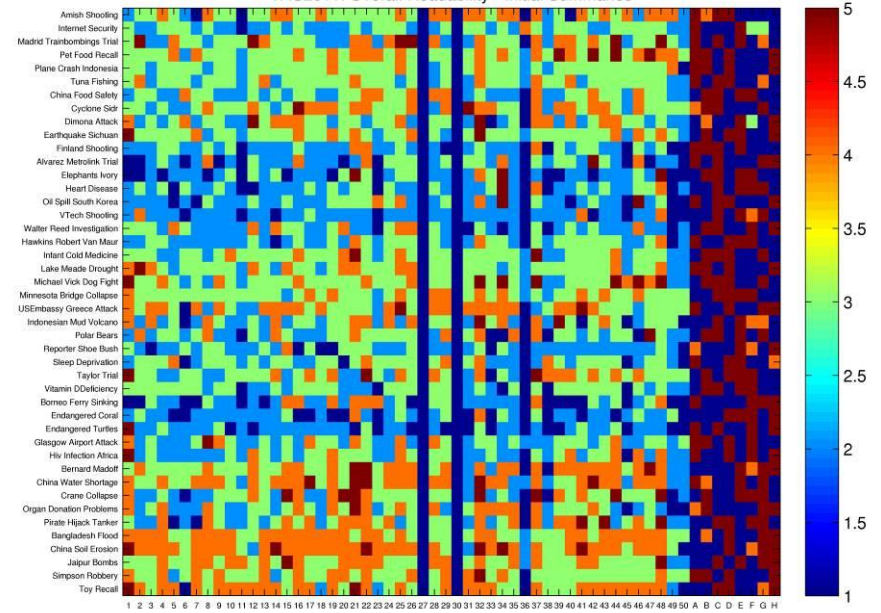
November 14, 2011

Motivation

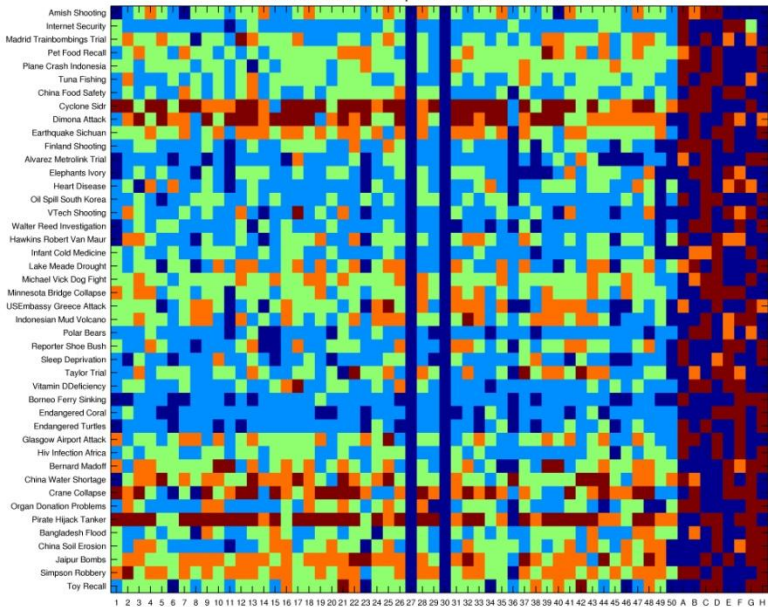
- Why do we propose using paired testing?
 - Paired testing can be more powerful than unpaired testing
 - Evidence suggests that document difficulty varies greatly, and paired testing accounts for this
 - Paired testing helps automatic metrics distinguish between humans and machines

Initial Summaries: The A Set

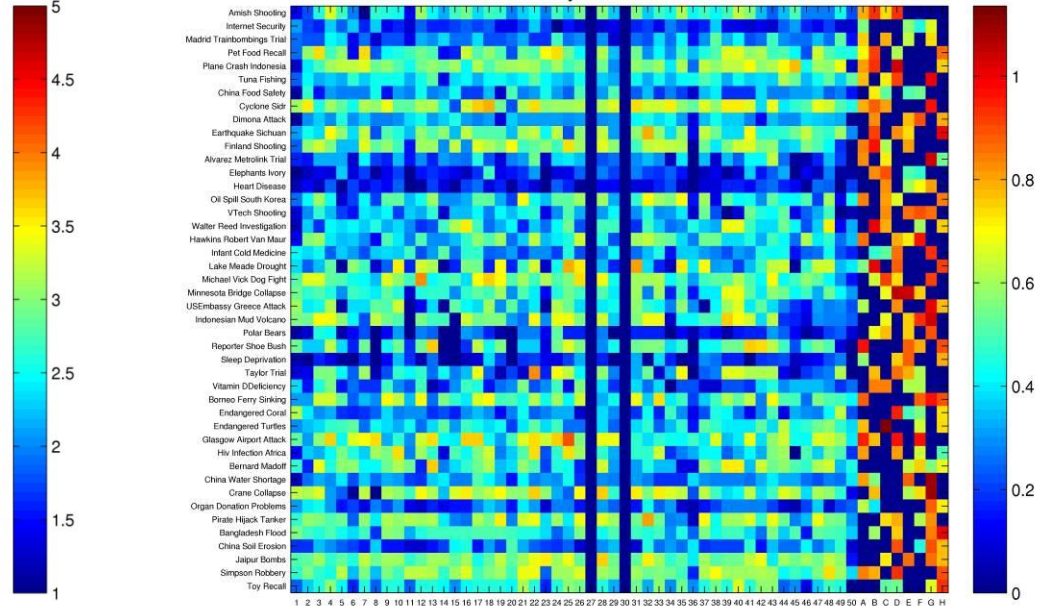
TAC2011: Overall Readability - Initial Summaries



TAC2011: Overall Responsiveness - Initial Summaries

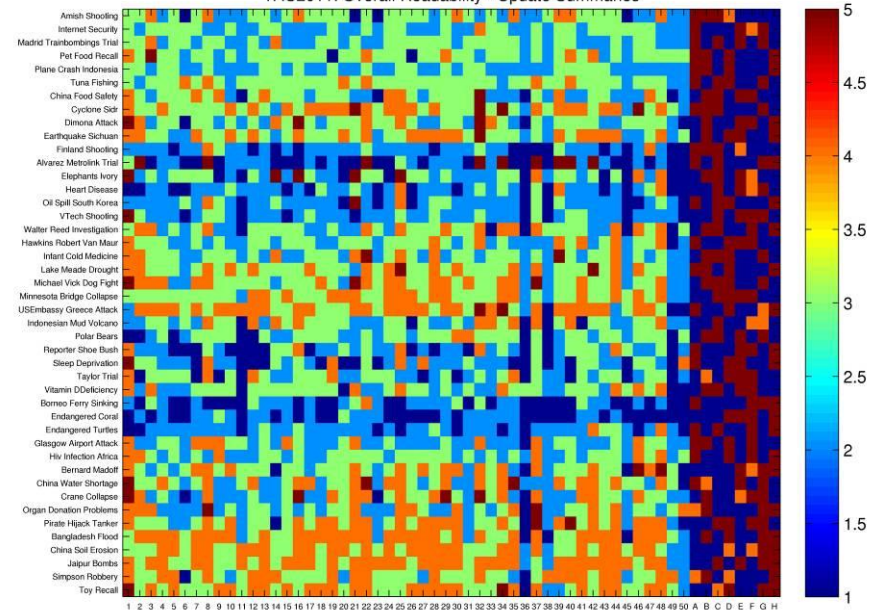


TAC2011: Pyramid - Initial Summaries

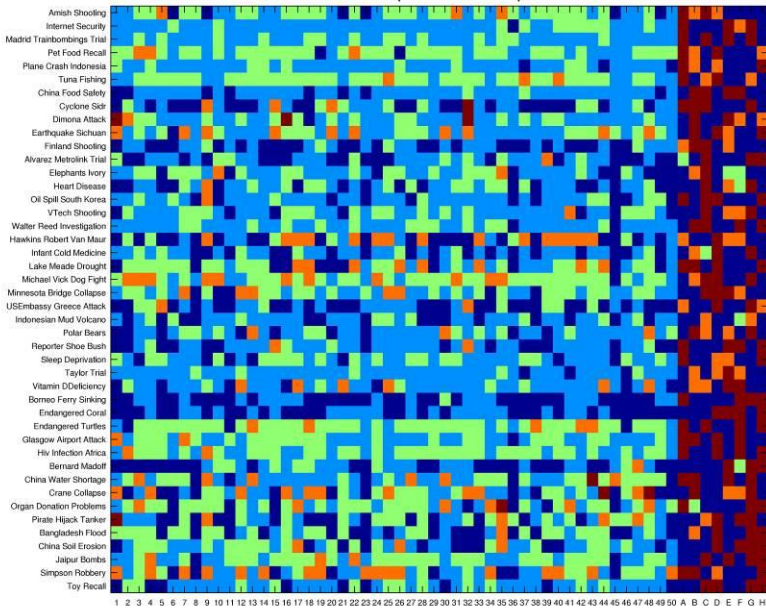


Update Summaries: The B Set

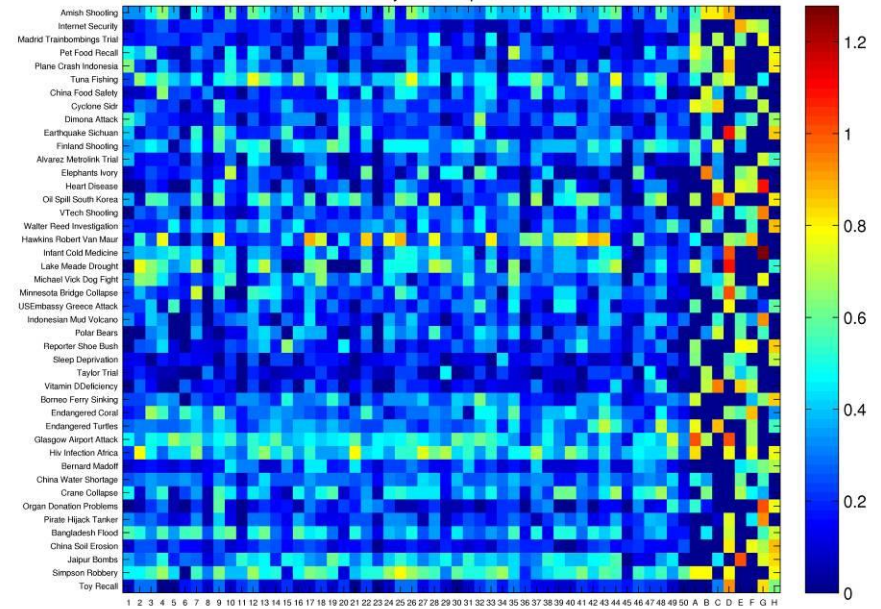
TAC2011: Overall Readability - Update Summaries



TAC2011: Overall Responsiveness - Update Summaries



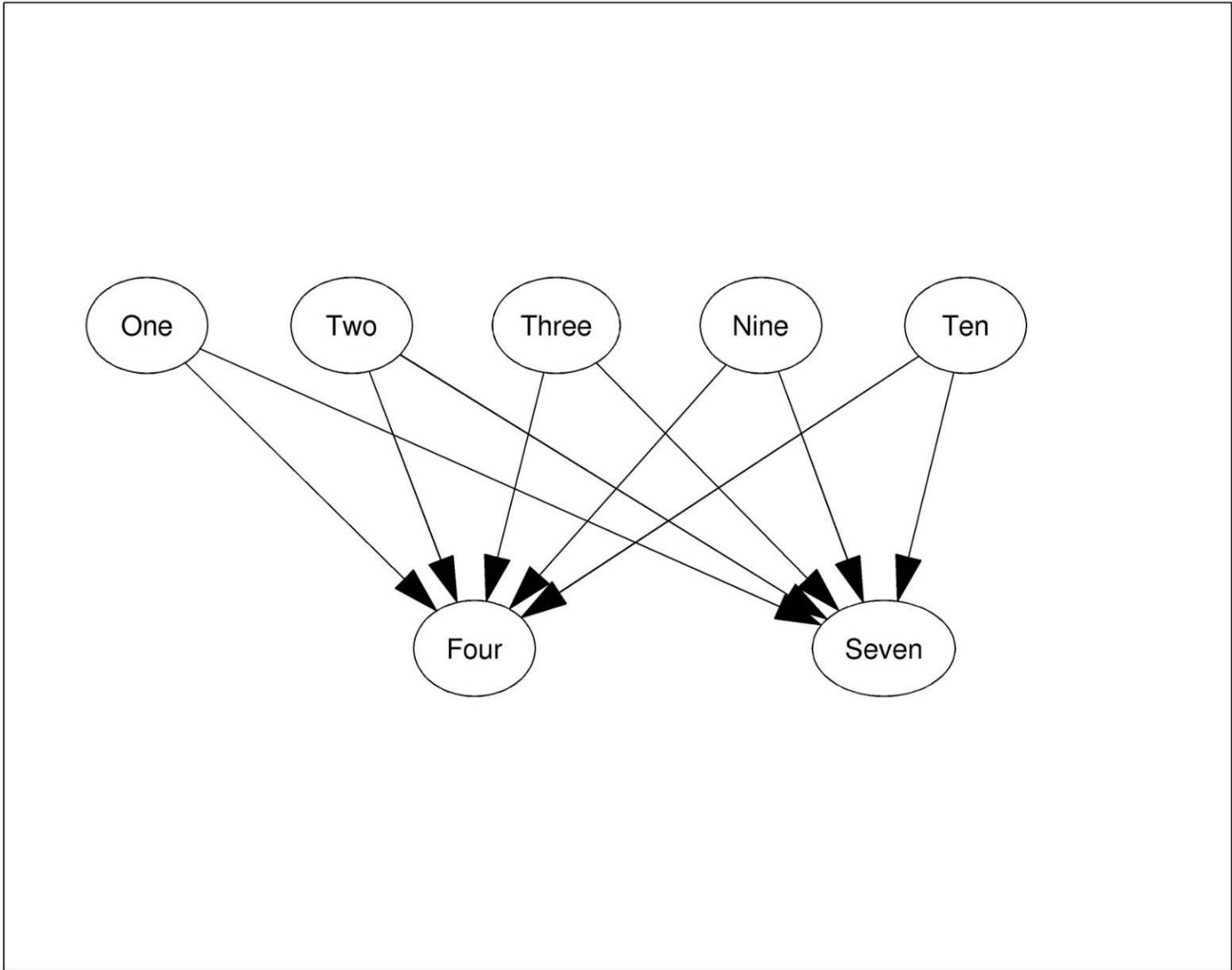
TAC2011: Pyramid - Update Summaries



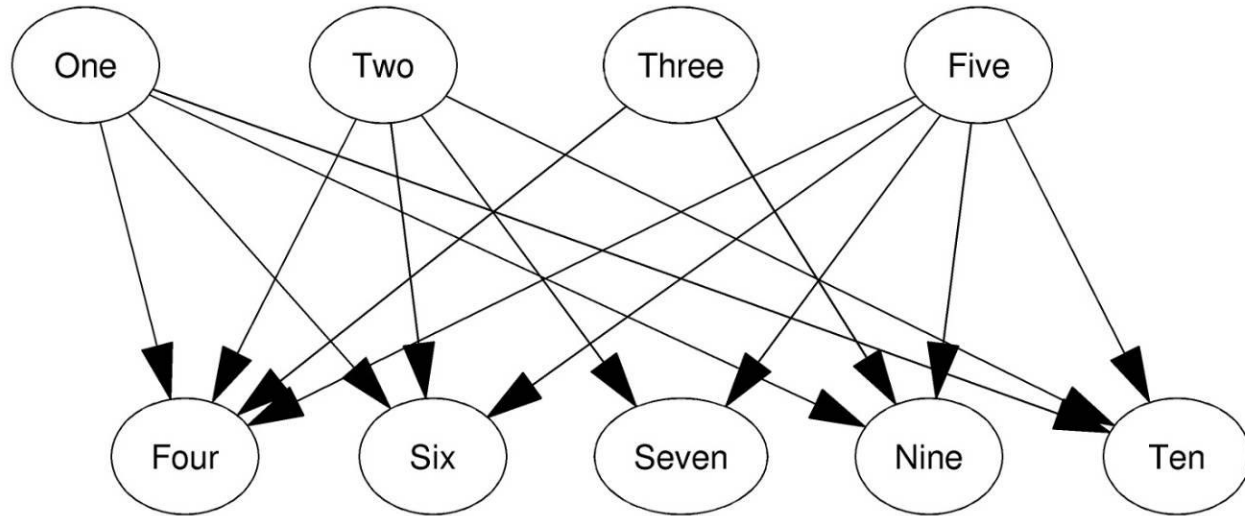
MultiLingual Task

- Use paired testing to evaluate system performance on individual languages
- Draw a directed graph displaying an edge (i,j) only when system i significantly outperformed system j
- Paired test was the non-parametric Wilcoxon signed-rank test
- Used Bonferroni adjustment to account for multiple comparisons

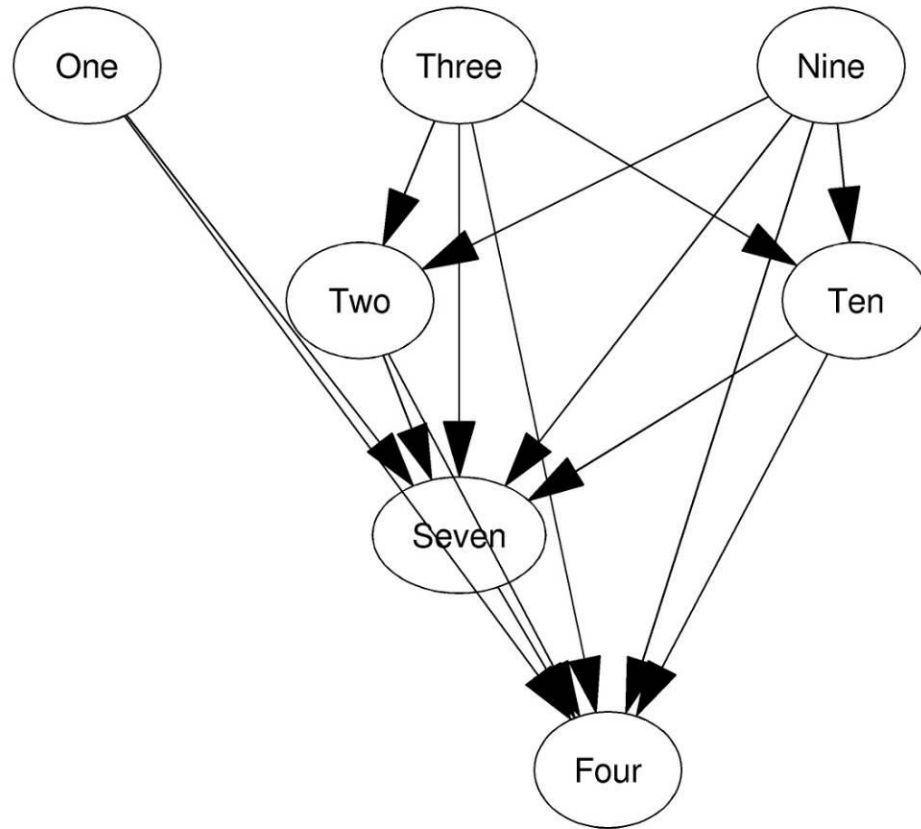
greek



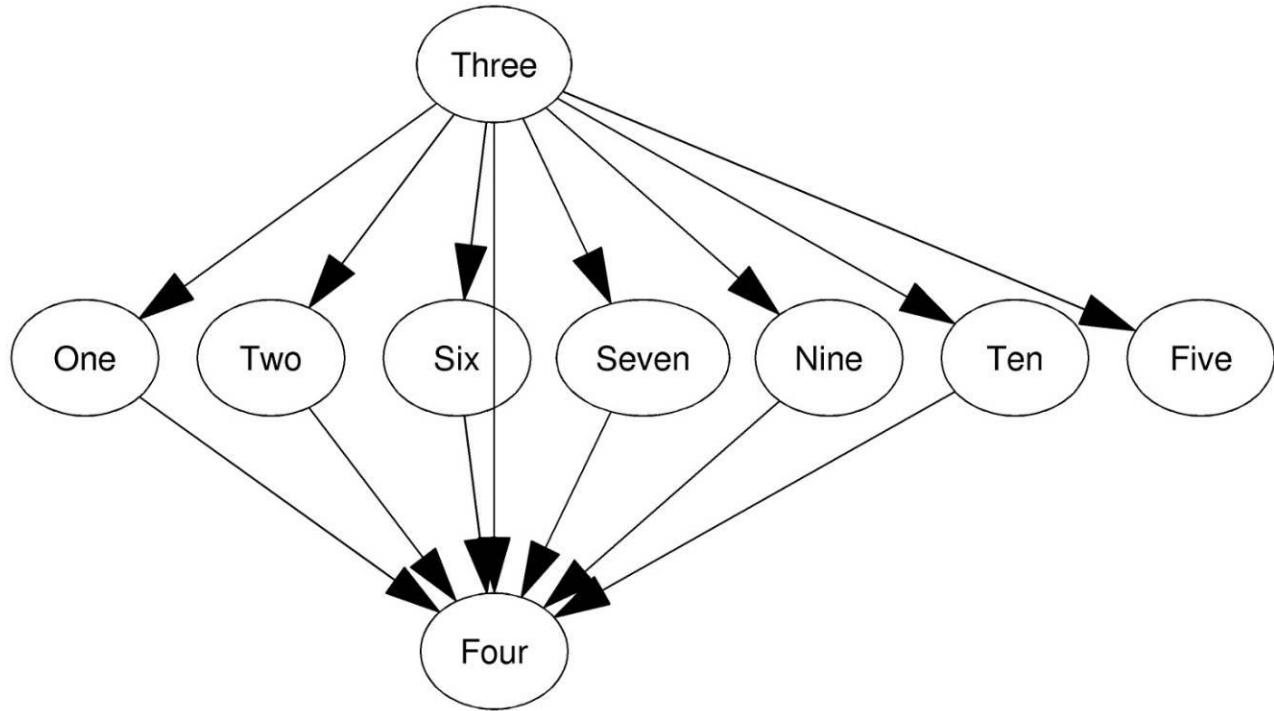
hindi



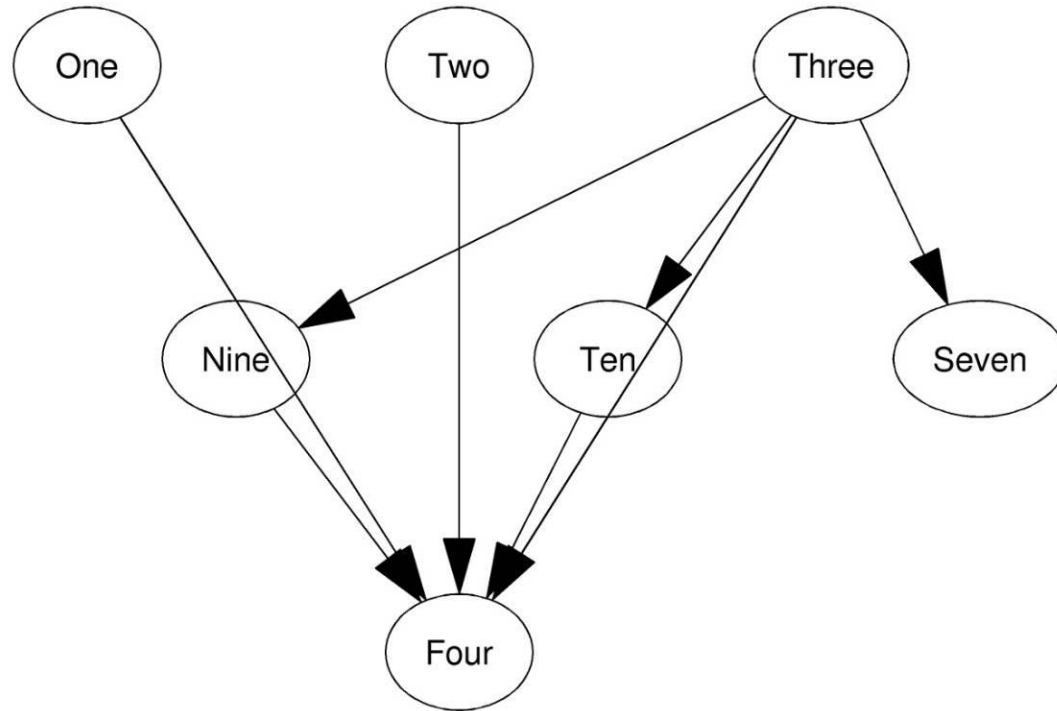
czech



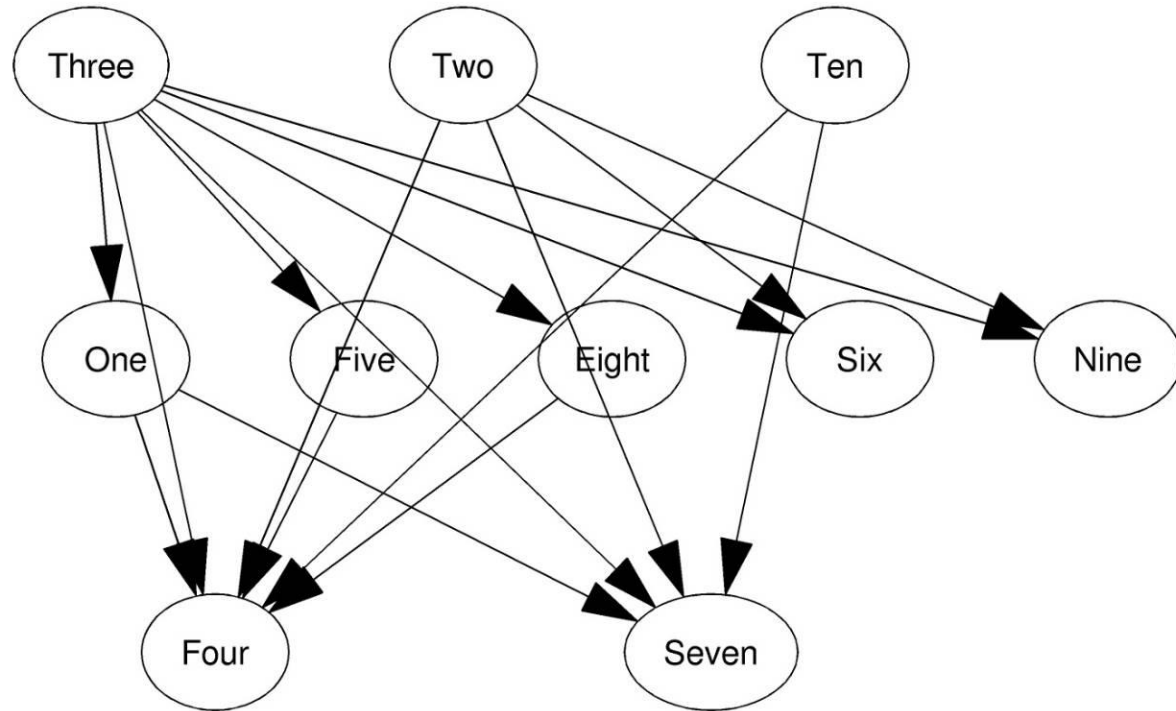
french



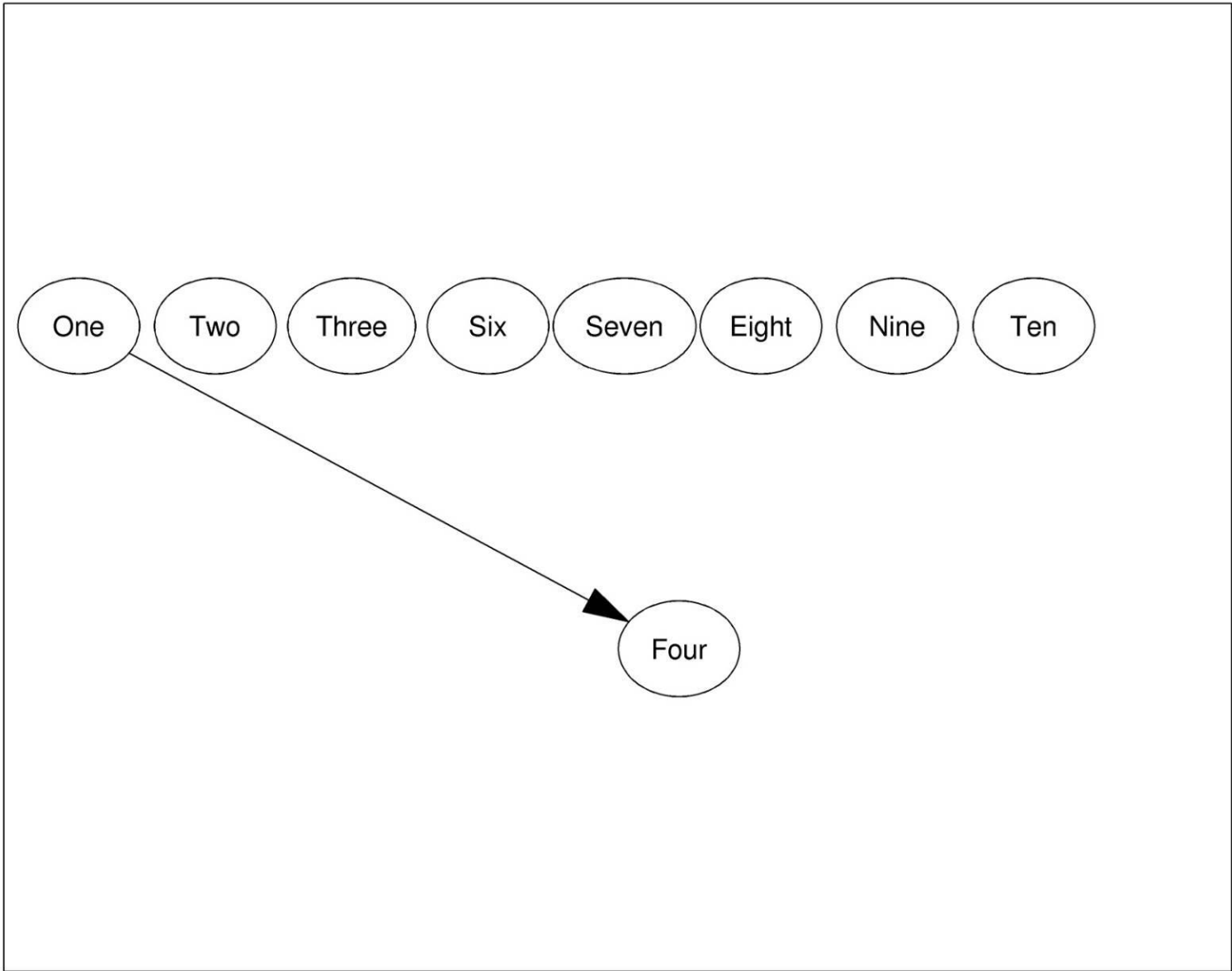
hebrew



english



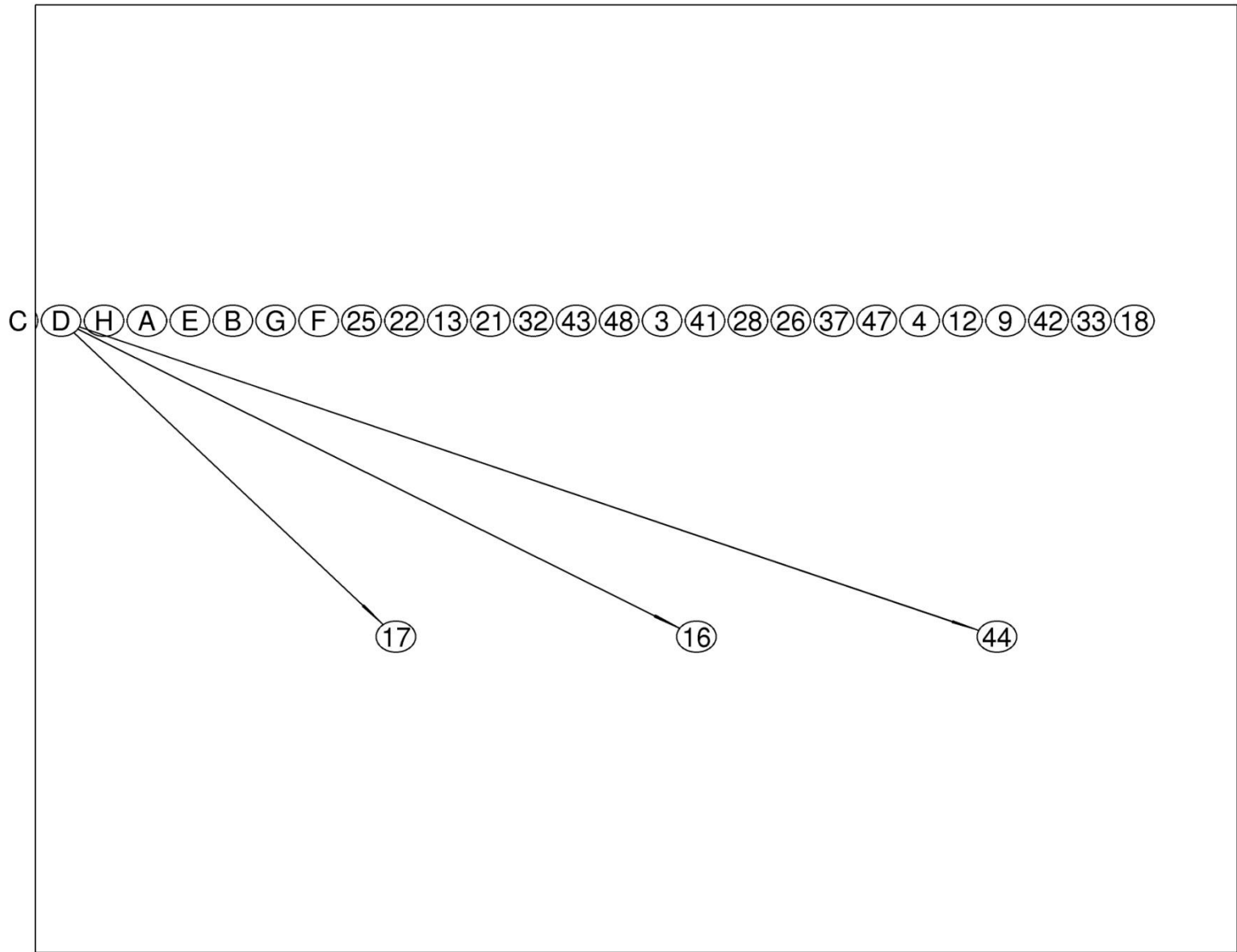
arabic



Guided Summarization

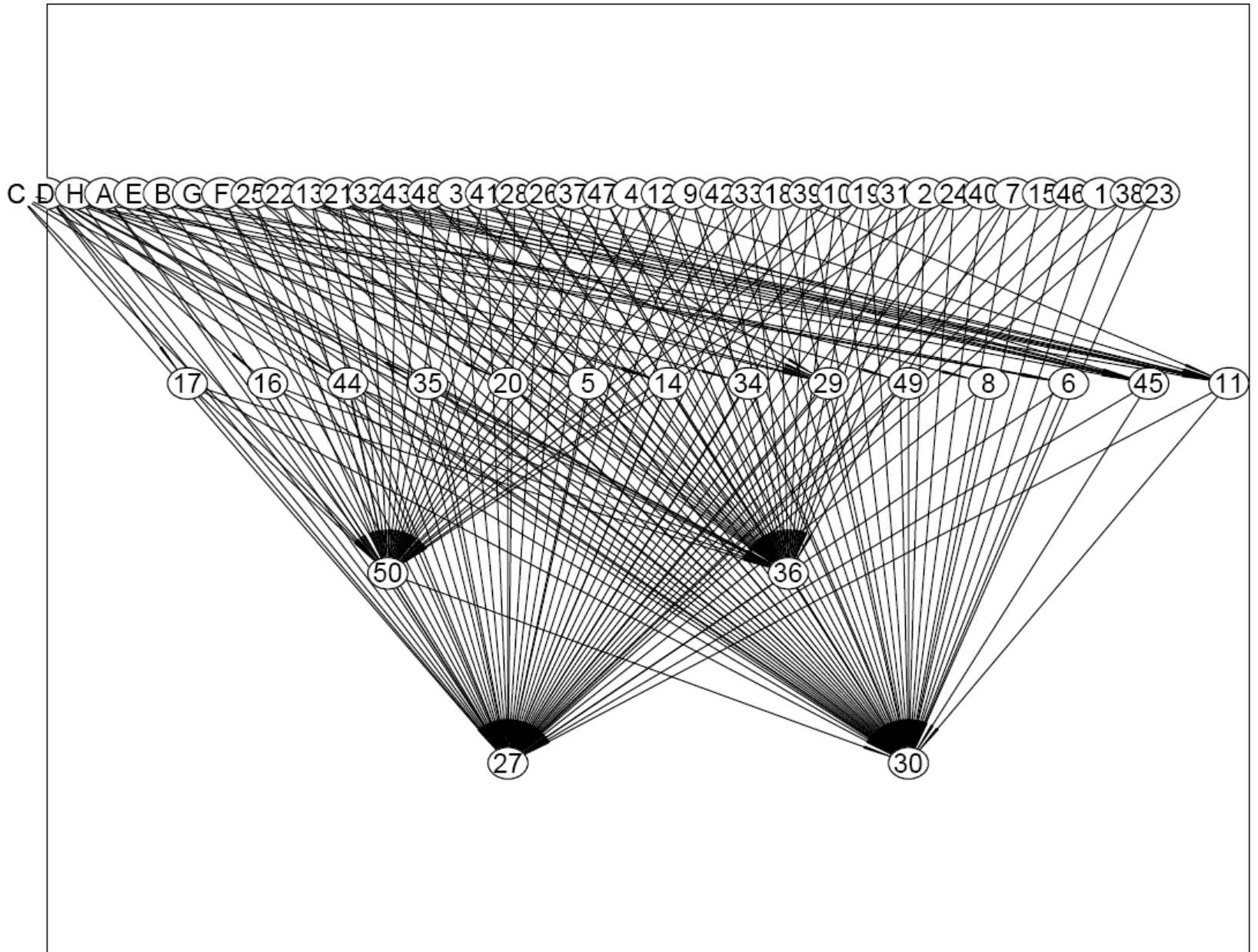
- Apply the same idea to systems in the Guided Summarization task
- Systems 1-50
- Humans A-H
- Wilcoxon signed-rank test
- Bonferroni Correction

Guided Summarizers Compared



The Top 30 Summarization Systems

Guided Summarizers Compared



Future Work

- More closely examine the AESOP results
 - Re-calculate the tables of discriminative significant differences using paired testing
- Look at TAC and DUC data from past years to see if there are trends.
 - E.g. Does it get more crowded at the top for a new task or a repeated task?

References

- Peter Rankel, John M. Conroy, Eric V. Slud, and Dianne P. O'Leary, “Ranking Human and Machine Summarization Systems,” [Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing \(EMNLP 2011\)](#), Edinburgh, UK, July 27-31, 2011. Association for Computational Linguistics (ACL) ([link to pdf](#))