

TECHNOLOGY-ASSISTED REVIEW IN E-DISCOVERY CAN BE
MORE EFFECTIVE AND MORE EFFICIENT
THAN EXHAUSTIVE MANUAL REVIEW

By Maura R. Grossman^{*} & Gordon V. Cormack^{†**}

Cite as: Maura R. Grossman & Gordon V. Cormack, *Technology-Assisted Review in E-Discovery Can Be More Effective and More Efficient Than Exhaustive Manual Review*, XVII RICH. J.L. & TECH. 11 (2011), <http://jolt.richmond.edu/v17i3/article11.pdf>.

^{*} Maura R. Grossman is counsel at Wachtell, Lipton, Rosen & Katz. She is co-chair of the E-Discovery Working Group advising the New York State Unified Court System, and a member of the Discovery Subcommittee of the Attorney Advisory Group to the Judicial Improvements Committee of the U.S. District Court for the Southern District of New York. Ms. Grossman is a coordinator of the Legal Track of the National Institute of Standards and Technology's Text Retrieval Conference ("TREC"), and an adjunct faculty member at Rutgers School of Law–Newark and Pace Law School. Ms. Grossman holds a J.D. from Georgetown University Law Center, and an M.A. and Ph.D. in Clinical/School Psychology from Adelphi University. The views expressed herein are solely those of the Author and should not be attributed to her firm or its clients.

[†] Gordon V. Cormack is a Professor at the David R. Cheriton School of Computer Science, and co-director of the Information Retrieval Group, at the University of Waterloo. He is a coordinator of the TREC Legal Track, and Program Committee member of TREC at large. Professor Cormack is the co-author of *Information Retrieval: Implementing and Evaluating Search Engines* (MIT Press, 2010), as well as more than 100 scholarly articles. Professor Cormack holds a B.Sc., M.Sc., and Ph.D. in Computer Science from the University of Manitoba.

^{**} The Authors would like to thank Ellen Voorhees and Ian Soboroff at NIST for providing access to the raw TREC 2009 data. The Authors gratefully acknowledge the helpful comments received from Hon. John M. Facciola (D.D.C.), Hon. Paul W. Grimm (D. Md.), and Hon. Andrew J. Peck (S.D.N.Y.) on an earlier draft of this paper.

ABSTRACT

E-discovery processes that use automated tools to prioritize and select documents for review are typically regarded as potential cost-savers – but inferior alternatives – to exhaustive manual review, in which a cadre of reviewers assesses every document for responsiveness to a production request, and for privilege. This Article offers evidence that such technology-assisted processes, while indeed more efficient, can also yield results superior to those of exhaustive manual review, as measured by recall and precision, as well as F_1 , a summary measure combining both recall and precision. The evidence derives from an analysis of data collected from the TREC 2009 Legal Track Interactive Task, and shows that, at TREC 2009, technology-assisted review processes enabled two participating teams to achieve results superior to those that could have been achieved through a manual review of the entire document collection by the official TREC assessors.

I. INTRODUCTION

[1] *The Sedona Conference Best Practices Commentary on the Use of Search and Information Retrieval Methods in E-Discovery* cautions that:

[T]here appears to be a myth that manual review by humans of large amounts of information is as accurate and complete as possible – perhaps even perfect – and constitutes the gold standard by which all searches should be measured. Even assuming that the profession had the time and resources to continue to conduct manual review of massive sets of electronic data sets (which it does not), the relative efficacy of that approach versus utilizing newly developed automated methods of review remains very much open to debate.¹

While the word *myth* suggests disbelief, literature on the subject contains little scientific evidence to support or refute the notion that automated methods, while improving on the efficiency of manual review, yield inferior results.² This Article presents evidence supporting the position that a technology-assisted process, in which humans examine only a small fraction of the document collection, can yield higher recall and/or precision than an exhaustive manual review process, in which humans code and examine the entire document collection.

[2] A *technology-assisted review process* involves the interplay of humans and computers to identify the documents in a collection that are responsive to a production request, or to identify those documents that should be withheld on the basis of privilege.³ A human examines and

¹ The Sedona Conference, *The Sedona Conference Best Practices Commentary on the Use of Search and Information Retrieval Methods in E-Discovery*, 8 SEDONA CONF. J. 189, 199 (2007) [hereinafter *Sedona Search Commentary*].

² *Id.* at 194 (“The comparative efficacy of the results of manual review versus the results of alternative forms of automated methods of review remains very much an open matter of debate.”).

³ See Douglas W. Oard et al., *Evaluation of information retrieval for E-discovery*, 18:4 ARTIFICIAL INTELLIGENCE & LAW 347, 365 (2010) (“In some cases . . . the end user will interact directly with the system, specifying the query, reviewing results, modifying the

codes only those documents the computer identifies – a tiny fraction of the entire collection.⁴ Using the results of this human review, the computer codes the remaining documents in the collection for responsiveness (or privilege).⁵ A technology-assisted review process may involve, in whole or in part, the use of one or more approaches including, but not limited to, keyword search, Boolean search, conceptual search, clustering, machine learning, relevance ranking, and sampling.⁶ In contrast, *exhaustive manual review* requires one or more humans to examine each and every document in the collection, and to code them as responsive (or privileged) or not.⁷

[3] Relevant literature suggests that manual review is far from perfect.⁸ Moreover, recent results from the Text Retrieval Conference (“TREC”), sponsored by the National Institute of Standards and Technology (“NIST”), show that technology-assisted processes can achieve high levels of recall and precision.⁹ By analyzing data collected

query, and so on. In other cases, the end user’s interaction with the system will be more indirect. . . .”).

⁴ See *Sedona Search Commentary supra* note 1, at 209.

⁵ See Maura R. Grossman & Terry Sweeney, *What Lawyers Need to Know About Search Tools*, THE NAT’L L.J. (Aug. 23, 2010), available at <http://www.law.com/jsp/lawtechnologynews/PubArticleLTN.jsp?id=1202470952987&slreturn=1&hbxlogin=1> (“‘machine learning tools,’ use ‘seed sets’ of documents previously identified as responsive or unresponsive to rank the remaining documents from most to least likely to be relevant, or to classify the documents as responsive or nonresponsive.”).

⁶ See, e.g., *Sedona Search Commentary, supra* note 1, at 217–23; CORNELIS JOOST VAN RIJSBERGEN, INFORMATION RETRIEVAL 74-85 (2d ed. 1979). The specific technologies employed in the processes that are the subjects of this study are detailed *infra* Parts III.A. – III.B.

⁷ See, e.g., Herbert L. Roitblat et al., *Document Categorization in Legal Electronic Discovery: Computer Classification vs. Manual Review*, 61 J. AM. SOC’Y. FOR INFO. SCI. AND TECH. 70, 70 (2010).

⁸ See, e.g., *Sedona Search Commentary, supra* note 1.

⁹ Bruce Hedin et al., *Overview of the TREC 2009 Legal Track*, in NIST SPECIAL PUBLICATION: SP 500-278, THE EIGHTEENTH TEXT RETRIEVAL CONFERENCE (TREC 2009) PROCEEDINGS 16 & tbl.5 (2009), available at <http://trec->

during the course of the TREC 2009 Legal Track Interactive Task,¹⁰ the Authors demonstrate that the levels of performance achieved by two technology-assisted processes exceed those that would have been achieved by the official TREC assessors – law students and lawyers employed by professional document-review companies – had they conducted a manual review of the entire document collection.

[4] Part II of this Article describes document review and production in the context of civil litigation, defines commonly used terms in the field of information retrieval, and provides an overview of recent studies. Part III details the TREC 2009 Legal Track Interactive Task, including the H5 and Waterloo efforts, as well as the TREC process for assessment and gold-standard creation. Part IV uses statistical inference to compare the recall, precision, and F_1 scores that H5 and Waterloo achieved to those the TREC assessors would have achieved had they reviewed all of the documents in the collection. Part V presents a qualitative analysis of the nature of manual review errors. Parts VI, VII, and VIII, respectively, discuss the results, limitations, and conclusions associated with this study. Ultimately, this Article addresses a fundamental uncertainty that arises in determining what is reasonable and proportional: Is it true that if a human examines every document from a particular source, that human will, as nearly as possible, correctly identify all and only the documents that should be produced? That is, does exhaustive manual review guarantee that production will be as complete and correct as possible? Or can technology-assisted review, in which a human examines only a fraction of the documents, do better?

II. CONTEXT

[5] Under Federal Rule of Civil Procedure 26(g)(1) (“Rule 26(g)(1)”), an attorney of record must certify “to the best of [his or her] knowledge,

legal.umiacs.umd.edu/LegalOverview09.pdf; *see also* Douglas W. Oard et al., *Overview of the TREC 2008 Legal Track*, in NIST SPECIAL PUBLICATION: SP 500-277, THE SEVENTEENTH TEXT RETRIEVAL CONFERENCE (TREC 2008) PROCEEDINGS 8 (2008), available at <http://trec.nist.gov/pubs/trec17/papers/LEGAL.OVERVIEW08.pdf>.

¹⁰ *See* Hedin et al., *supra* note 9, at 2.

information, and belief formed after a reasonable inquiry,” that every discovery request, response, or objection is

consistent with [the Federal Rules of Civil Procedure] . . . not interposed for any improper purpose, such as to harass, cause unnecessary delay, or needlessly increase the cost of litigation[, and is] neither unreasonable nor unduly burdensome or expensive, considering the needs of the case, prior discovery in the case, the amount in controversy, and the importance of the issues at stake in the action.¹¹

Similarly, Federal Rule of Civil Procedure 26(b)(2)(C)(iii) (“Rule 26(b)(2)(C)(iii)”) requires a court to limit discovery when it determines that “the burden or expense of the proposed discovery outweighs its likely benefit, considering the needs of the case, the amount in controversy, the parties’ resources, the importance of the issues at stake in the action, and the importance of the discovery in resolving the issues.”¹² Thus, Rules 26(g)(1) and 26(b)(2)(C)(iii) require that discovery requests and responses be *proportional*.¹³ However, Federal Rule of Civil Procedure 37(a)(4) (“Rule 37(a)(4)”) provides that “an evasive or incomplete disclosure, answer or response must be treated as a failure to disclose, answer, or respond[.]” and therefore requires that discovery responses be *complete*.¹⁴ Together, Rules 26(g)(1), 26(b)(2)(C)(iii), and 37(a)(4) reflect the tension – between completeness on one hand, and burden and cost on the other – that exists in all electronic discovery (“e-discovery”) processes.¹⁵ In

¹¹ FED. R. CIV. P. 26(g)(1).

¹² FED. R. CIV. P. 26(b)(2)(C)(iii).

¹³ The Sedona Conference, *The Sedona Conference Commentary on Proportionality in Electronic Discovery*, 11 SEDONA CONF. J. 289, 294 (2010) [hereinafter *Sedona Proportionality Commentary*].

¹⁴ FED. R. CIV. P. 37(a)(4).

¹⁵ Typically, a responding party will not only seek to produce *all* responsive documents, but to identify *only* the responsive documents, in order to guard against overproduction or waiver of privilege. *See, e.g., Mt. Hawley Ins. Co. v. Felman Prod., Inc.*, 271 F.R.D. 125, 136 (S.D.W. Va. 2010) (finding that plaintiff’s over-production of documents by more than 30% was a factor in waiver of privilege).

assessing what is reasonable and proportional with respect to e-discovery, parties and courts must balance these competing considerations.¹⁶

[6] One of the greatest challenges facing legal stakeholders is determining whether or not the cost and burden of identifying and producing electronically stored information (“ESI”) is commensurate with its importance in resolving the issues in dispute.¹⁷ In current practice, the problem of identifying responsive (or privileged) ESI, once it has been collected, is almost always addressed, at least in part, by a manual review process, the cost of which dominates the e-discovery process.¹⁸ A natural question to ask, then, is whether this manual review process is the most effective and efficient one for identifying and producing the ESI most likely to resolve a dispute.

A. Information Retrieval

[7] The task of finding all, and only, the documents that meet “some requirement” is one of information retrieval (“IR”), a subject of scholarly

¹⁶ See *Harkabi v. Sandisk Corp.*, No. 08 Civ. 8203 (WHP), 2010 WL 3377338, at *1 (S.D.N.Y. Aug. 23, 2010) (“Electronic discovery requires litigants to scour disparate data storage mediums and formats for potentially relevant documents. That undertaking involves dueling considerations: thoroughness and cost.”).

¹⁷ See *id.* at *8 (“Integral to a court’s inherent power is the power to ensure that the game is worth the candle—that commercial litigation makes economic sense. Electronic discovery in this case has already put that principle in jeopardy.”); *Hopson v. Mayor of Balt.*, 232 F.R.D. 228, 232 (D. Md. 2005) (“This case vividly illustrates one of the most challenging aspects of discovery of electronically stored information—how properly to conduct Rule 34 discovery within a reasonable pretrial schedule, while concomitantly insuring that requesting parties receive appropriate discovery, and that producing parties are not subjected to production timetables that create unreasonable burden, expense, and risk of waiver of attorney-client privilege and work product protection”). See generally *Sedona Proportionality Commentary*, *supra* note 13.

¹⁸ Marisa Peacock, *The True Cost of eDiscovery*, CMSWiRE, <http://www.cmswire.com/cms/enterprise-cms/the-true-cost-of-ediscovery-006060.php> (2009) (citing *Sedona Search Commentary*, *supra* note 1, at 192); Ashish Prasad et al., *Cutting to the “Document Review” Chase: Managing a Document Review in Litigation and Investigations*, 18 BUS. LAW TODAY, 2, Nov.–Dec. 2008.

research for at least a century.¹⁹ In IR terms, “some requirement” is referred to as an *information need*, and *relevance* is the property of whether or not a particular document meets the information need.²⁰ For e-discovery, the information need is typically specified by a production request (or by the rules governing privilege), and the definition of relevance follows.²¹ Cast in IR terms, the objective of review in e-discovery is to identify as many *relevant* documents as possible, while simultaneously identifying as few *nonrelevant* documents as possible.²² The fraction of relevant documents identified during a review is known as *recall*, while the fraction of identified documents that are relevant is known as *precision*.²³ That is, *recall* is a measure of completeness, while *precision* is a measure of accuracy, or correctness.²⁴

[8] The notion of *relevance*, although central to information science, and the subject of much philosophical and scientific investigation, remains elusive.²⁵ While it is easy enough to write a document describing an

¹⁹ The concepts and terminology outlined in Part II.A may be found in many information retrieval textbooks. For a historical perspective, see GERARD SALTON & MICHAEL J. MCGILL, *INTRODUCTION TO MODERN INFORMATION RETRIEVAL* (1983); VAN RIJSBERGEN, *supra* note 6. For a more modern treatment, see STEFAN BÜTTCHER ET AL., *INFORMATION RETRIEVAL: IMPLEMENTING AND EVALUATING SEARCH ENGINES* 33–75 (2010).

²⁰ See BÜTTCHER ET AL., *supra* note 19, at 5-6, 8.

²¹ See Hedin et al., *supra* note 9, at 1.

²² See VAN RIJSBERGEN, *supra* note 6, at 4.

²³ See David C. Blair & M. E. Maron, *An Evaluation of Retrieval Effectiveness for a Full-Text Document-Retrieval System*, 28 *COMMUN. ACM* 289, 290 (1985) (“Recall measures how well a system retrieves *all* the relevant documents; and Precision, how well the system retrieves *only* the relevant documents.”); VAN RIJSBERGEN, *supra* note 6, at 112-13.

²⁴ See VAN RIJSBERGEN, *supra* note 6, at 113.

²⁵ See Tefko Saracevic, *Relevance: A Review of the Literature and a Framework for Thinking on the Notion in Information Science. Part II: Nature and Manifestations of Relevance*, 58 *J. AM. SOC’Y FOR INFO. SCI. & TECH.* 1915 (2007); Tefko Saracevic, *Relevance: A Review of the Literature and a Framework for Thinking on the Notion in*

information need and hence relevance, determining the relevance of any particular document requires human interpretation.²⁶ It is well established that human assessors will disagree in a substantial number of cases as to whether a document is relevant, regardless of the information need or the assessors' expertise and diligence.²⁷

[9] A review resulting in higher recall and higher precision than another review is more nearly complete and correct, and therefore superior,²⁸ while a review with lower recall and lower precision is inferior.²⁹ If one result has higher recall while the other has higher precision, it is not immediately obvious which should be considered superior. To calculate a review's effectiveness, researchers often employ F_1 – the harmonic mean of recall and precision³⁰ – a commonly used summary measure that rewards results achieving both high recall and high precision, while penalizing those that have either low recall or low precision.³¹ The value of F_1 is always intermediate between recall and precision, but is generally closer to the lesser of the two.³² For example, a result with 40% recall and 60% precision has $F_1 = 48\%$. Following

Information Science. Part III: Behavior and Effects of Relevance, 58:13 J. AM. SOC'Y FOR INFO. SCI. & TECH. 2126 (2007).

²⁶ See Peter Bailey et al., *Relevance Assessment: Are Judges Exchangeable and Does It Matter?*, in SIGIR '08 PROCEEDINGS OF THE 31ST ANNUAL INTERNATIONAL ACM SIGIR CONFERENCE ON RESEARCH AND DEVELOPMENT IN INFORMATION RETRIEVAL 667 (2008); see also VAN RIJSBERGEN, *supra* note 6, at 112.

²⁷ See Bailey et al., *supra* note 26, at § 4.3.

²⁸ See Blair & Maron, *supra* note 23.

²⁹ See *id.*

$$^{30} F_1 = \frac{2}{\frac{1}{\text{recall}} + \frac{1}{\text{precision}}}.$$

³¹ See BÜTTCHER ET AL., *supra* note 19, at 68.

³² See *id.*

TREC, this Article reports recall and precision, along with F_1 as a summary measure of overall review effectiveness.³³

B. Assessor Overlap

[10] The level of agreement between independent assessors may be quantified by *overlap* – also known as the *Jaccard index* – the number of documents identified as relevant by two independent assessors, divided by the number identified as relevant by either or both assessors.³⁴ For example, suppose assessor A identifies documents {W,X,Y,Z} as relevant, while assessor B identifies documents {V,W,X}. Both assessors have identified two documents {W,X} as relevant, while either or both have identified five documents {V,W,X,Y,Z} as relevant. So the overlap is 2/5, or forty percent. Informally, overlap of less than fifty percent indicates that the assessors disagree on whether or not a document is relevant more often than when they agree that a document is relevant.³⁵

[11] In her study, *Variations in Relevance Judgments and the Measurement of Retrieval Effectiveness*, Ellen Voorhees measured overlap between primary, secondary, and tertiary reviewers who each made 14,968 assessments of relevance for 13,435 documents,³⁶ with respect to 49

³³ See Hedin et al., *supra* note 9, at 3.

³⁴ Ellen M. Voorhees, *Variations in Relevance Judgments and the Measurement of Retrieval Effectiveness*, 36 INFO. PROCESSING & MGMT 697, 700 (2000), available at http://www.cs.cornell.edu/courses/cs430/2006fa/cache/Trec_8.pdf (“Overlap is defined as the size of the intersection of the relevant document sets divided by the size of the union of the relevant document sets.”); see CHRISTOPHER D. MANNING ET AL., AN INTRODUCTION TO INFORMATION RETRIEVAL 61 (2009) (draft), available at nlp.stanford.edu/IR-book/pdf/irbookonlinereading.pdf; see also Raimundo Real & Juan M. Vargas, *The Probabilistic Basis of Jaccard’s Index of Similarity*, 45 SYSTEMATIC BIOLOGY 380, 381 (1996).

³⁵ See Ellen M. Voorhees, *The Philosophy of Information Retrieval Evaluation*, in EVALUATION OF CROSS-LANGUAGE INFORMATION RETRIEVAL SYSTEMS SECOND WORKSHOP OF THE CROSS-LANGUAGE EVALUATION FORUM, CLEF 2001 DARMSTADT, GERMANY, SEPTEMBER 3-4, 2001 REVISED PAPERS 355, 364 (Carol Peters et al. eds., 2002).

³⁶ E-mail from Ellen M. Voorhees to Gordon V. Cormack (Jul. 31, 2019 14:34 EDT) (on file with authors). The numbers in the text are derived from the file,

information needs (or “topics,” in TREC parlance), in connection with Ad Hoc Task of the Fourth Text Retrieval Conference (“TREC 4”).³⁷ As illustrated in Table 1, the overlap between primary and secondary assessors was 42.1%;³⁸ the overlap between primary and tertiary assessors was 49.4%;³⁹ and the overlap between secondary and tertiary assessors was 42.6%.⁴⁰

[12] Perhaps due to the assessors’ expertise,⁴¹ Voorhees’ overlap results are among the highest reported for pairs of human assessors. Her findings demonstrate that assessors disagree at least as often as they agree that a document is relevant.⁴² Voorhees concluded:

The scores for the [secondary and tertiary] judgments imply a practical upper bound on retrieval system performance is 65% precision at 65% recall since that is the level at which humans agree with one another.⁴³

“threeWayJudgments,” attached to Voorhees’ e-mail. Some of the documents were assessed for relevance to more than one topic.

³⁷ Voorhees, *supra* note 34, at 708; *see also* Donna Harman, *Overview of the Fourth Text REtrieval Conference (TREC-4)*, in NIST SPECIAL PUBLICATION 500-236: THE FOURTH TEXT RETRIEVAL CONFERENCE (TREC-4) 2 (2004), *available at* http://trec.nist.gov/pubs/trec4/t4_proceedings.html (follow the first link under “PAPERS”).

³⁸ *See infra* Table 1; *see also* Voorhees, *supra* note 34, at 701 tbl.1.

³⁹ *See infra* Table 1; *see also* Voorhees, *supra* note 34, at 701 tbl.1.

⁴⁰ *See infra* Table 1; *see also* Voorhees, *supra* note 34, at 701 tbl.1.

⁴¹ All assessors were professional information retrieval experts. Voorhees, *supra* note 34, at 701.

⁴² *See id.*

⁴³ *Id.*

[13] It is not widely accepted that these findings apply to e-discovery.⁴⁴ This “legal exceptionalism” appears to arise from common assumptions within the legal community:

1. that the information need (responsiveness or privilege) is more precisely defined for e-discovery than for classical information retrieval;⁴⁵
2. that lawyers are better able to assess relevance and privilege than the non-lawyers typically employed for information retrieval tasks;⁴⁶ and
3. that the most defensible way to ensure that a production is accurate is to have a lawyer examine each and every document.⁴⁷

⁴⁴ See *Sedona Search Commentary*, *supra* note 1 (noting the widespread perception that manual review is nearly perfect). If that perception were correct, manual reviewers would have close to 100% overlap, contrary to Voorhees’ findings. Voorhees, *supra* note 34, at 701 tbl.1.

⁴⁵ Oard et al., *supra* note 3, at 362 (“It is important to recognize that the notion of relevance that is operative in E-discovery is, naturally, somewhat more focused than what has been studied in information seeking behavior studies generally . . .”).

⁴⁶ Cf. Alejandra P. Perez, *Assigning Non-Attorneys to First-Line Document Reviews Requires Safeguards*, THE E-DISCOVERY 4-1-1 (LeClairRyan), Jan. 2011, at 1, available at <http://marketing.leclairryan.com/files/Uploads/Documents/the-e-discovery-4-1-1-01-21-2011.pdf> (opining that non-attorney document reviewers typically require additional training, particularly regarding the legal concept of privilege).

⁴⁷ See *Sedona Search Commentary*, *supra* note 1, at 203 (“Some litigators continue to primarily rely upon manual review of information as part of their review process. Principal rationales [include] . . . the perception that there is a lack of scientific validity of search technologies necessary to defend against a court challenge”); see also Thomas E. Stevens & Wayne C. Matus, *A ‘Comparative Advantage’ To Cut E-Discovery Costs*, NAT’L L.J. (Sept. 4, 2008), <http://www.law.com/jsp/nlj/PubArticleNLJ.jsp?id=1202424251053> (describing a “general reluctance by counsel to rely on anything but what they perceive to be the most defensible positions in electronic discovery, even if those solutions do not hold up any sort of honest analysis of cost or quality”).

Assumptions (1) and (2) are amenable to scientific evaluation, as is the overarching question of whether technology-assisted review can improve upon exhaustive manual review. Assumption (3) – a legal opinion – should be informed by scientific evaluation of the first two assumptions.

| Assessment | Primary | Secondary | Tertiary |
|------------|---------|-----------|----------|
| Primary | 100% | | |
| Secondary | 42.1% | 100% | |
| Tertiary | 49.4% | 42.6% | 100% |

Table 1: Overlap in relevance assessments by primary, secondary, and tertiary assessors for the TREC 4 Ad Hoc Task.⁴⁸

[14] Recently, Herbert Roitblat, Anne Kershaw, and Patrick Oot studied the level of agreement among review teams using data produced to the Department of Justice (“DOJ”) in response to a Second Request that stemmed from MCI’s acquisition of Verizon.⁴⁹ In their study, two independent teams of professional assessors, Teams A and B, reviewed a random sample of 5,000 documents.⁵⁰ Roitblat and his colleagues reported the level of agreement and disagreement between the original production, Team A, and Team B, as a contingency matrix,⁵¹ from which the Authors calculated overlap, as shown in Table 2.⁵² The overlap between Team A and the original production was 16.3%,⁵³ the overlap between Team B and the original production was 15.8%;⁵⁴ and the overlap between Teams A and B was 28.1%.⁵⁵ These and other studies of overlap

⁴⁸ Voorhees, *supra* note 34, at 701 tbl.1.

⁴⁹ See Roitblat et al., *supra* note 7, at 73.

⁵⁰ See *id.* at 73-74.

⁵¹ *Id.* at 74 tbl.1.

⁵² See *infra* Table 2.

⁵³ *Id.*

⁵⁴ *Id.*

⁵⁵ *Id.*

indicate that relevance is not a concept that can be applied consistently by independent assessors, even if the information need is specified by a production request and the assessors are lawyers.⁵⁶

| Assessment | Production | Team A | Team B |
|------------|------------|--------|--------|
| Production | 100% | | |
| Team A | 16.3% | 100% | |
| Team B | 15.8% | 28.1% | 100% |

Table 2: Overlap in relevance assessments between original production in a Second Request, and two subsequent manual reviews.⁵⁷

C. Assessor Accuracy

[15] Measurements of overlap provide little information regarding the accuracy of particular assessors because there is no “gold standard” against which to compare them.⁵⁸ One way to resolve this problem is to deem one assessor’s judgments correct by definition, and to use those judgments as the gold standard for the purpose of evaluating the other assessor(s).⁵⁹

[16] In the Voorhees study, the primary assessor composed the information need specification for each topic.⁶⁰ It may therefore be reasonable to take the primary assessor’s coding decisions to be the gold standard. In the Roitblat, Kershaw, and Oot study, a senior attorney familiar with the case adjudicated all instances of disagreement between Teams A and B.⁶¹ Although Roitblat and his colleagues sought to

⁵⁶ See Roitblat et al., *supra* note 7, at 73; Voorhees, *supra* note 34.

⁵⁷ The Authors derived the information in Table 2 from the Roitblat, Kershaw, and Oot study. Roitblat et al., *supra* note 7, at 74; *see supra* para. 13.

⁵⁸ Roitblat et al., *supra* note 7, at 77.

⁵⁹ See Voorhees, *supra* note 34, at 700.

⁶⁰ *Id.*

⁶¹ Roitblat et al., *supra* note 7, at 74.

measure agreement,⁶² it may be reasonable to use their “adjudicated results” as the gold standard. These adjudicated results deemed the senior attorney’s opinion correct in cases where Teams A and B disagreed, and deemed the consensus correct in cases where Teams A and B agreed.⁶³ Assuming these gold standards, Table 3 shows the effectiveness of the various assessors in terms of recall, precision, and F_1 .⁶⁴ Note that recall ranges from 52.8% to 83.6%, while precision ranges from 55.5% to 81.9%, and F_1 ranges from 64.0% to 70.4%.⁶⁵ All in all, these results appear to be reasonable, but hardly perfect. Can technology-assisted review improve on them?

D. Technology-Assisted Review Accuracy

[17] In addition to the two manual review groups, Roitblat, Kershaw, and Oot had two service providers (Teams C and D) use technology-assisted review processes to classify each document in the dataset as

⁶² *Id.* at 72 (“Formally, the present study is intended to examine the hypothesis: *The rate of agreement between two independent reviewers of the same documents will be equal to or less than the agreement between a computer-aided system and the original review.*”).

⁶³ *Id.* at 74.

The 1,487 documents on which Teams A and B disagreed were submitted to a senior Verizon litigator (P. Oot), who adjudicated between the two teams, again without knowledge of the specific decisions made about each document during the first review. This reviewer had knowledge of the specifics of the matter under review, but had not participated in the original review. This authoritative reviewer was charged with determining which of the two teams had made the correct decision.

Id.

⁶⁴ *See infra* Table 3. Recall and precision for the secondary and tertiary assessors, using the primary assessor as the gold standard, are provided by Voorhees, *supra* note 34, at 701 tbl.2; recall and precision for Teams A and B, using the adjudicated results as the gold standard, were derived from Roitblat et al., *supra* note 7, at 74 tbl.1; F_1 was calculated from recall and precision using the formula at *supra* note 30.

⁶⁵ *See infra* Table 3.

relevant or not.⁶⁶ Unfortunately, the adjudicated results described in Part II.C. were made available to one of the two service providers, and therefore, cannot be used as a gold standard to evaluate the accuracy of the providers' efforts.⁶⁷

| Study | Review | Recall | Precision | F_1 |
|-----------------|-----------|--------|-----------|-------|
| Voorhees | Secondary | 52.8% | 81.3% | 64.0% |
| Voorhees | Tertiary | 61.8% | 81.9% | 70.4% |
| Roitblat et al. | Team A | 77.1% | 60.9% | 68.0% |
| Roitblat et al. | Team B | 83.6% | 55.5% | 66.7% |

Table 3: Recall, precision, and F_1 of manual assessments in studies by Voorhees, and Roitblat et al. Voorhees evaluated secondary and tertiary assessors with respect to a primary assessor, who was deemed correct. The Authors computed recall, precision, and F_1 from the results reported by Roitblat et al. for Teams A and B, using their adjudicated results as the gold standard.⁶⁸

[18] Instead, Roitblat and his colleagues reported recall, precision, and F_1 using, as an alternate gold standard, the set of documents originally produced to, and accepted by, the DOJ.⁶⁹ There is little reason to believe that this original production, and hence the alternate gold standard, was perfect.⁷⁰ The first two rows of Table 4 show the recall and precision of manual review Teams A and B when evaluated with respect to this

⁶⁶ Roitblat et al., *supra* note 7, at 74-75.

⁶⁷ *Id.* at 74 (“One of these systems based its classifications in part on the adjudicated results of Teams A and B, but without any knowledge of how those teams’ decisions were related to the decisions made by [the] original review team. As a result, it is not reasonable to compare the classifications of these two systems to the classifications of the two re-review teams, but it is reasonable to compare them to the classifications of the original review.”).

⁶⁸ Voorhees, *supra* note 34, at 701 tbl.2; Roitblat et al. *supra* note 7, at 74 tbl.1.

⁶⁹ Roitblat et al., *supra* note 7, at 74.

⁷⁰ *Id.* at 76 (“The use of precision and recall implies the availability of a stable ground truth against which to compare the assessments. Given the known variability of human judgments, we do not believe that we have a solid enough foundation to claim that we know which documents are truly relevant and which are not.”).

alternate gold standard.⁷¹ These results are much worse than those in Table 3.⁷² Team A achieved 48.8% recall and 19.7% precision, while Team B achieved 52.9% recall and 18.3% precision.⁷³ The corresponding F_1 scores were 28.1% and 27.2%, respectively – less than half of the F_1 scores achieved with respect to the gold standard derived using the senior attorney’s opinion.⁷⁴

[19] The recall and precision Roitblat, Kershaw, and Oot reported were computed using the original production as the gold standard, and are dramatically different from those shown in Table 3, which were computed using their adjudicated results as the gold standard.⁷⁵ Nevertheless, both sets of results appear to suggest the *relative* accuracy between Teams A and B: Team B has higher recall, while Team A has higher precision and higher F_1 , regardless of which gold standard is applied.⁷⁶

[20] The last two rows of Table 4 show the effectiveness of the technology-assisted reviews conducted by teams C and D, as reported by Roitblat, Kershaw, and Oot using the original production as the gold standard.⁷⁷ The results suggest that technology-assisted review Teams C and D achieved about the same recall as manual review Teams A and B, and somewhat better precision and F_1 .⁷⁸ However, due to the use of the alternate gold standard, the result is inconclusive.⁷⁹ Because the

⁷¹ See *id.* at 76 tbl.2; *infra* Table 4.

⁷² Compare *supra* Table 3, with *infra* Table 4.

⁷³ See *infra* Table 4; see also Roitblat et al., *supra* note 7, at 74-76.

⁷⁴ Compare *supra* Table 3, with *infra* Table 4.

⁷⁵ Compare *supra* Table 3, with *infra* Table 4. See generally Roitblat et al., *supra* note 7, at 76 tbl.2.

⁷⁶ See *supra* Table 3; *infra* Table 4; Roitblat et al., *supra* note 7, at 76 tbl.2.

⁷⁷ See *infra* Table 4; see also Roitblat et al., *supra* note 7, at 74-75.

⁷⁸ See *infra* Table 4.

⁷⁹ See Roitblat et al., *supra* note 7, at 76 (“The use of precision and recall implies the availability of a stable ground truth against which to compare the assessments. Given the

improvement from using technology-assisted review, as reported by Roitblat and his colleagues, is small compared to the difference between the results observed using the two different gold standards, it is difficult to determine whether the improvement represents a real difference in effectiveness as compared to manual review.

| Study | Review | Method | Recall | Precision | F_1 |
|-----------------|--------|-------------|--------|-----------|-------|
| Roitblat et al. | Team A | Manual | 48.8% | 19.7% | 28.1% |
| Roitblat et al. | Team B | Manual | 52.9% | 18.3% | 27.2% |
| Roitblat et al. | Team C | Tech. Asst. | 45.8% | 27.1% | 34.1% |
| Roitblat et al. | Team D | Tech. Asst. | 52.7% | 29.5% | 37.8% |

Table 4: Recall, precision, and F_1 of manual and technology-assisted review teams, evaluated with respect to the original production to the DOJ. The first two rows of this table differ from the last two rows of Table 3 only in the gold standard used for evaluation.⁸⁰

[21] In a heavily cited study by David C. Blair and M.E. Maron, skilled paralegal searchers were instructed to retrieve at least 75% of all documents relevant to 51 requests for information pertaining to a legal matter.⁸¹ For each request, the searchers composed keyword searches using an interactive search system, retrieving and printing documents for further review.⁸² This process was repeated until the searcher was satisfied that 75% of the relevant documents had been retrieved.⁸³ Although the searchers believed they had found 75% of the relevant documents, their average recall was only 20.0%.⁸⁴ Despite this low rate of

known variability of human judgments, we do not believe that we have a solid enough foundation to claim that we know which documents are truly relevant and which are not.”).

⁸⁰ *Id.* at 73-76.

⁸¹ *See* Blair & Maron, *supra* note 23, at 291.

⁸² *Id.*

⁸³ *Id.*

⁸⁴ *Id.* at 293; *see also* Maureen Dostert & Diane Kelly, *Users' Stopping Behaviors and Estimates of Recall*, in SIGIR '09 PROCEEDINGS OF THE 32ND ANNUAL INTERNATIONAL ACM SIGIR CONFERENCE ON RESEARCH AND DEVELOPMENT IN INFORMATION RETRIEVAL 820–21 (2009) (showing that most subjects in an interactive information

recall, the searchers achieved a high average precision of 79.0%.⁸⁵ From the published data,⁸⁶ the Authors calculated the average F_1 score to be 28.0% – remarkably similar to that observed by Roitblat and his colleagues for manual review.⁸⁷

[22] Blair and Maron argue that the searchers would have been unable to achieve higher recall even if they had known there were many relevant documents that were not retrieved.⁸⁸ Researcher Gerald Salton disagrees.⁸⁹ He claims that it would have been possible for the searchers to achieve higher recall at the expense of lower precision, either by broadening their queries or by taking advantage of the relevance ranking capability of the search system.⁹⁰

[23] Overall, the literature offers little reason to believe that manual review is perfect. But is it as complete and accurate as possible, or can it be improved upon by technology-assisted approaches invented since Blair and Maron's study?

[24] As previously noted, recent results from TREC suggest that technology-assisted approaches may indeed be able to improve on manual review.⁹¹ In the TREC 2008 Legal Track Interactive Task, H5, a San

retrieval experiment reported they had found about 51-60% of the relevant documents when, on average, recall was only 7%).

⁸⁵ See Blair & Maron, *supra* note 23, at 293.

⁸⁶ *Id.*

⁸⁷ See Roitblat et al., *supra* note 7 at 76.

⁸⁸ See Blair & Maron, *supra* note 23, at 295-96.

⁸⁹ See Gerard Salton, *Another Look at Automatic Text-Retrieval Systems*, 29:7 COMM'NS ACM 648, 650 (1986).

⁹⁰ *Id.* at 648-49.

⁹¹ See generally Hedin et al., *supra* note 9; Oard et al., *supra* note 9.

Francisco-based legal information retrieval firm,⁹² employed a user-modeling approach⁹³ to achieve recall, precision, and F_1 of 62.4%, 81.0%, and 70.5%, respectively, in response to a mock request to produce documents from a 6,910,192-document collection released under the tobacco Master Settlement Agreement.⁹⁴ In the course of this effort, H5 examined only 7,992 documents⁹⁵ – roughly 860 times fewer than the 6,910,192 it would have been necessary to examine in an exhaustive manual review. Yet the results compare favorably with those previously reported for manual review or keyword search, exceeding what Voorhees characterizes as a “practical upper bound” on what may be achieved, given uncertainties in assessment.⁹⁶

⁹² See *Contact Us*, H5, <http://www.h5.com/about/contact.php> (last visited Mar. 22, 2011); *Who We Are*, H5, http://www.h5.com/about/who_we_are.html (last visited Apr. 11, 2011).

⁹³ Christopher Hogan et al., *H5 at TREC 2008 Legal Interactive: User Modeling, Assessment & Measurement*, in NIST SPECIAL PUBLICATION: SP 500-277, THE SEVENTEENTH TEXT RETRIEVAL CONFERENCE (TREC 2008) PROCEEDINGS (2008), available at <http://trec.nist.gov/pubs/trec17/papers/h5.legal.rev.pdf> (last visited Mar. 23, 2011).

⁹⁴ Oard et al., *supra* note 9, at 30 tbl.15; see also *Complex Document Image Processing (CDIP)*, ILL. INST. TECH., <http://ir.iit.edu/projects/CDIP.html> (last visited Apr. 11, 2011); *Master Settlement Agreement*, NAT’L ASS’N ATTORNEYS GEN. (Nov. 1998), available at <http://www.naag.org/backpages/naag/tobacco/msa/msa-pdf/MSA%20with%20Sig%20Pages%20and%20Exhibits.pdf>; TREC 2008, *Complaint for Violation of the Federal Securities Laws, Mellon v. Echinoderm Cigarettes, Inc.*, (2008), available at <http://trec-legal.umiacs.umd.edu/topics/8I.pdf>.

⁹⁵ Hogan et al., *supra* note 92, at 8.

⁹⁶ Voorhees, *supra* note 34, at 701.

| Topic | Production Request |
|-------|---|
| 201 | All documents or communications that describe, discuss, refer to, report on, or relate to the Company's engagement in structured commodity transactions known as "prepay transactions." |
| 202 | All documents or communications that describe, discuss, refer to, report on, or relate to the Company's engagement in transactions that the Company characterized as compliant with FAS 140 (or its predecessor FAS 125). |
| 203 | All documents or communications that describe, discuss, refer to, report on, or relate to whether the Company had met, or could, would, or might meet its financial forecasts, models, projections, or plans at any time after January 1, 1999. |
| 204 | All documents or communications that describe, discuss, refer to, report on, or relate to any intentions, plans, efforts, or activities involving the alteration, destruction, retention, lack of retention, deletion, or shredding of documents or other evidence, whether in hard-copy or electronic form. |
| 205 | All documents or communications that describe, discuss, refer to, report on, or relate to energy schedules and bids, including but not limited to, estimates, forecasts, descriptions, characterizations, analyses, evaluations, projections, plans, and reports on the volume(s) or geographic location(s) of energy loads. |
| 206 | All documents or communications that describe, discuss, refer to, report on, or relate to any discussion(s), communication(s), or contact(s) with financial analyst(s), or with the firm(s) that employ them, regarding (i) the Company's financial condition, (ii) analysts' coverage of the Company and/or its financial condition, (iii) analysts' rating of the Company's stock, or (iv) the impact of an analyst's coverage of the Company on the business relationship between the Company and the firm that employs the analyst. |
| 207 | All documents or communications that describe, discuss, refer to, report on, or relate to fantasy football, gambling on football, and related activities, including but not limited to, football teams, football players, football games, football statistics, and football performance. |

Table 5: Mock production requests ("topics") composed for the TREC 2009 Legal Track Interactive Task.⁹⁷

⁹⁷ TREC 2009, *Complaint, Grumby v. Volteron Corp.*, 14 (2009) available at http://trec-legal.umiacs.umd.edu/LT09_Complaint_J_final.pdf; see also Hedin et al., *supra* note 9, at 5-6.

[25] One of the Authors was inspired to try to reproduce these results at TREC 2009 using an entirely different approach: statistical active learning, originally developed for e-mail spam filtering.⁹⁸ At the same time, H5 reprised its approach for TREC 2009.⁹⁹ The TREC 2009 Legal Track Interactive Task used the same design as TREC 2008, but employed a different complaint¹⁰⁰ and seven new mock requests to produce documents (see Table 5) from a new collection of 836,165 e-mail messages and attachments captured from Enron at the time of its collapse.¹⁰¹ Each participating team was permitted to request as many topics as they wished, however, due to resource constraints, the most topics that any team was assigned was four of the seven.¹⁰²

⁹⁸ See generally Gordon V. Cormack & Mona Mojdeh, *Machine Learning for Information Retrieval: TREC 2009 Web, Relevance Feedback and Legal Tracks*, in NIST SPECIAL PUBLICATION: SP 500-278, THE EIGHTEENTH TEXT RETRIEVAL CONFERENCE (TREC 2009) PROCEEDINGS (2009), available at <http://trec.nist.gov/pubs/trec18/papers/uwaterloo-cormack.WEB.RF.LEGAL.pdf>.

⁹⁹ Hedin et al., *supra* note 9, at 6.

¹⁰⁰ See generally TREC 2009, *Complaint*, *supra* note 97.

¹⁰¹ Hedin et al., *supra* note 9, at 4; see *Information Released in Enron Investigation*, FED. ENERGY REG. COMM'N, <http://www.ferc.gov/industries/electric/indus-act/wec/enron/info-release.asp> (last visited Apr. 11, 2011) [hereinafter FEREC]; E-mail from Bruce Hedin to Gordon V. Cormack (Aug. 31, 2009 20:33 EDT) (on file with authors) (“I have attached full list of the 836,165 document-level IDs . . .”). The collection is available at *Practice Topic and Assessments for TREC 2010 Legal Learning Task*, U. WATERLOO, <http://plg1.uwaterloo.ca/~gvcormac/treclegal09/> (follow “The TREC 2009 dataset”) (last visited Apr. 18, 2011).

¹⁰² Hedin et al., *supra* note 9, at 7; E-mail from Bruce Hedin to Gordon V. Cormack & Maura R. Grossman (Mar. 24, 2011 02:46 EDT) (on file with authors).

| Team | Topic | Reviewed | Produced | Recall | Precision | F_1 |
|----------|----------|----------|----------|--------|-----------|-------|
| Waterloo | 201 | 6,145 | 2,154 | 77.8% | 91.2% | 84.0% |
| Waterloo | 202 | 12,646 | 8,746 | 67.3% | 88.4% | 76.4% |
| Waterloo | 203 | 4,369 | 2,719 | 86.5% | 69.2% | 76.9% |
| H5 | 204 | 20,000 | 2,994 | 76.2% | 84.4% | 80.1% |
| Waterloo | 207 | 34,446 | 23,252 | 76.1% | 90.7% | 82.8% |
| | Average: | 15,521 | 7,973 | 76.7% | 84.7% | 80.0% |

Table 6: Effectiveness of H5 and Waterloo submissions to the TREC 2009 Legal Track Interactive Task.¹⁰³

[26] Together, H5 and Waterloo produced documents for five distinct TREC 2009 topics;¹⁰⁴ the results of their efforts are summarized in Table 6. The five efforts employed technology-assisted processes, with the number of manually reviewed documents for each topic ranging from 4,369 to 34,446¹⁰⁵ (or 0.5% to 4.1% of the collection). That is, the total human effort for the technology-assisted processes – measured by the number of documents reviewed – was between 0.5% and 4.1% of that which would have been necessary for an exhaustive manual review of all 836,165 documents in the collection.¹⁰⁶ The number of documents produced for each topic ranged from 2,154 to 23,252¹⁰⁷ (or 0.3% to 2.8% of the collection; about half the number of documents reviewed). Over the five efforts, the average recall and precision were 76.7% and 84.7%,

¹⁰³ See *infra*, para. 25.

¹⁰⁴ See Hedin et al., *supra* note 9, at 7.

¹⁰⁵ Cormack & Mojdeh, *supra* note 98, at 6 tbl.2 (showing that Waterloo reviewed between 4,369 documents (for Topic 203) and 34,446 documents (for Topic 207); see E-mail from Dan Brassil to Maura R. Grossman (Dec. 17, 2010 15:21 EST) (on file with authors) (“[H5] sampled and reviewed 20,000 documents”).

¹⁰⁶ See sources cited *supra* note 101.

¹⁰⁷ NIST Special Publication 500-277: The Seventeenth Text REtrieval Conference Proceedings (TREC 2008) http://trec.nist.gov/pubs/trec17/t17_proceedings.html Appendix: Per Topic Scores: TREC 2009 Legal Track, Interactive Task, 3 tbl.4, 4 tbl.8, 5 tbl.12, 6 tbl.16, 9 tbl.26 <http://trec.nist.gov/pubs/trec18/appendices/app09int2.pdf>.

respectively; no recall was lower than 67.3%, and no precision was lower than 69.2%,¹⁰⁸ placing all five efforts above what Voorhees characterized as a “practical upper bound” on what may be achieved, given uncertainties in assessment.¹⁰⁹

[27] Although it appears that the TREC results are better than those previously reported in the literature, either for manual or technology-assisted review, they do not include any direct comparison between manual and technology-assisted review.¹¹⁰ To draw any firm conclusion that one is superior to the other, one must compare manual and technology-assisted review efforts using the same information needs, the same dataset, and the same evaluation standard.¹¹¹ The Roitblat, Kershaw, and Oot study is the only peer-reviewed study known to the Authors suggesting that technology-assisted review *may be* superior to manual review – if only in terms of precision, and only by a small amount – based on a common information need, a common dataset, and a common gold standard, albeit one of questionable accuracy.¹¹²

[28] This Article shows conclusively that the H5 and Waterloo efforts *are* superior to manual reviews conducted contemporaneously by TREC assessors, using the same topics, the same datasets, and the same gold standard. The manual reviews considered for this Article were the “First-Pass Assessments” undertaken at the request of the TREC coordinators for

¹⁰⁸ See Hedin et al, *supra* note 9, at 17.

¹⁰⁹ Voorhees, *supra* note 34, at 701.

¹¹⁰ See e.g., Oard et al., *supra* note 9, at 1-2.

¹¹¹ See Voorhees, *supra* note 35, at 356 (“The [Cranfield] experimental design called for the same set of documents and same set of information needs to be used for each [search method], and for the use of both precision and recall to evaluate the effectiveness of the search.”).

¹¹² See Roitblat et al., *supra* note 7, at 76 (“The use of precision and recall implies the availability of a stable ground truth against which to compare the assessments. Given the known variability of human judgments, we do not believe that we have a solid enough foundation to claim that we know which documents are truly relevant and which are not.”).

the purpose of evaluating the participating teams' submissions.¹¹³ In comparing the manual and technology-assisted reviews, the Authors used exactly the same adjudicated gold standard as TREC.¹¹⁴

III. TREC Legal Track Interactive Task

[29] TREC is an annual event hosted by NIST, with the following objectives:

- to encourage research in information retrieval based on large test collections;
- to increase communication among industry, academia, and government by creating an open forum for the exchange of research ideas;
- to speed the transfer of technology from research labs into commercial products by demonstrating substantial improvements in retrieval methodologies on real-world problems; and
- to increase the availability of appropriate evaluation techniques for use by industry and academia, including development of new evaluation techniques more applicable to current systems.¹¹⁵

Since its inception in 2006,¹¹⁶ the TREC Legal Track has had the goal “to develop search technology that meets the needs of lawyers to engage in effective discovery in digital document collections.”¹¹⁷

¹¹³ Hedin et al., *supra* note 9, at 3 (describing the “First-Pass Assessment” process).

¹¹⁴ *See id.* at 3-4.

¹¹⁵ Text REtrieval Conference (TREC), *Overview*, NAT'L INST. STANDARDS & TECH., <http://trec.nist.gov/overview.html> (last updated Aug. 10, 2010).

¹¹⁶ *See* Jason R. Baron, *The TREC Legal Track: Origins and Reflections on the First Year*, 8 SEDONA CONF. J. 251, 253 (2007); *see also* Jason R. Baron et al., *TREC-2006 Legal Track Overview*, in NIST SPECIAL PUBLICATION: SP 500-272, THE FIFTEENTH TEXT RETRIEVAL CONFERENCE (TREC 2006) PROCEEDINGS 1-2 (2006), *available at* <http://trec.nist.gov/pubs/trec15/papers/LEGAL06.OVERVIEW.pdf>.

[30] Within the TREC Legal Track, the Interactive Task simulates the process of review of a large population of documents for responsiveness to one or more discovery requests in a civil litigation.¹¹⁸ In 2008, the first year of the Interactive Task,¹¹⁹ the population of documents used was the “Illinois Institute of Technology Complex Document Information Processing Test Collection, version 1.0” (“IIT CDIP”),¹²⁰ consisting of about seven million documents that were released in connection with various lawsuits filed against certain U.S. tobacco companies and affiliated research institutes.¹²¹ A mock complaint and three associated requests for production (or topics) were composed for the purposes of the Interactive Task.¹²² Participating teams were required to produce the responsive documents for one or more of the three requests.¹²³

[31] The population of documents used for TREC 2009 consisted of e-mail messages and attachments that Enron produced in response to requests by FERC.¹²⁴ A mock complaint and seven associated requests for production were composed for the purposes of TREC 2009.¹²⁵ Participating teams requested as many topics as they desired to undertake, but time and cost constraints limited the number of topics that any team was assigned to a maximum of four.¹²⁶

¹¹⁷ Text Retrieval Conference (TREC), *TREC Tracks*, NAT’L INST. STANDARDS & TECH., <http://trec.nist.gov/tracks.html> (last updated Feb. 24, 2011).

¹¹⁸ See Oard et al., *supra* note 9, at 20.

¹¹⁹ See *id.* at 2.

¹²⁰ *Id.* at 3; see *Complex Document Image Processing (CDIP)*, *supra* note 94.

¹²¹ See Oard et al., *supra* note 9, at 3; *Complex Document Image Processing (CDIP)*, *supra* note 93.

¹²² See Oard et al., *supra* note 9 at 3, 24.

¹²³ *Id.* at 24.

¹²⁴ See Hedin et al., *supra* note 9, at 4; see also FERC, *supra* note 101.

¹²⁵ See Hedin et al., *supra* note 9, at 5-6.

¹²⁶ See *id.* at 7 tbl.1.

[32] Aside from the document collections, the mock complaints, and the production requests, the conduct of the 2008 and 2009 Interactive Tasks was identical.¹²⁷ Participating teams were given the document collection, the complaint, and the production requests several weeks before production was due.¹²⁸ Teams were allowed to use any combination of technology and human input; the exact combination differed from team to team.¹²⁹ However, the size of the document population, along with time and cost constraints, rendered it infeasible for any team to conduct an exhaustive review of every document.¹³⁰ To the Authors' knowledge, no team examined more than a small percentage of the document population; H5 and Waterloo, in particular, used various combinations of computer search, knowledge engineering, machine learning, and sampling to select documents for manual review.¹³¹

[33] To aid the teams in their efforts, as well as to render an authoritative interpretation of responsiveness (or relevance, within the context of TREC), a volunteer *Topic Authority* ("TA") – a senior attorney familiar with the subject matter – was assigned for each topic.¹³² The TA played three critical roles:

- to consult with the participating teams to clarify the notion of relevance, in a manner chosen by the teams;

¹²⁷ See *id.* at 1-2.

¹²⁸ See Text Retrieval Conference (TREC), *TREC-2008 Legal Track Interactive Task: Guidelines*, 8, 17 (2008), trec-legal.umiacs.umd.edu/2008InteractiveGuidelines.pdf [hereinafter *TREC-2008 Guidelines*]; see also E-mail from Dan Brassil to Maura R. Grossman, *supra* note 105.

¹²⁹ *TREC-2008 Guidelines*, *supra* note 128, at 4, 7; see also E-mail from Bruce Hedin to Gordon V. Cormack (Apr. 07, 2011 00:56 EDT) (confirming that teams were permitted to use any combination of technology and human input).

¹³⁰ See *TREC-2008 Legal Track Interactive Task: Guidelines*, *supra* note 128, at 8.

¹³¹ See Hogan et al., *supra* note 9, at 5; Cormack & Mojdeh, *supra* note 98, at 6.

¹³² See Hedin et al., *supra* note 9, at 2.

- to prepare a set of written guidelines used by the human reviewers to evaluate, after the fact, the relevance of documents produced by the teams; and
- to act as a final arbiter of relevance in the adjudication process.¹³³

[34] The TREC coordinators evaluated the various participant efforts using estimates of recall, precision, and F_1 based on a two-pass human assessment process.¹³⁴ In the first pass, human reviewers assessed a stratified sample of about 7,000 documents for relevance.¹³⁵ For some topics (Topics 201, 202, 205, and 206), the reviewers were primarily volunteer law students supervised by the TREC coordinators; for others (Topics 203, 204, and 207), the reviewers were lawyers employed and supervised by professional document-review companies, who volunteered their services.¹³⁶

[35] The TREC coordinators released the first-pass assessments to participating teams, which were invited to appeal relevance determinations with which they disagreed.¹³⁷ For each topic, the TA adjudicated the appeals, and the TA's opinion was deemed to be correct and final.¹³⁸ The gold standard of relevance for the documents in each sample was therefore:

- The same as the first-pass assessment, for any document that participants did not appeal; or

¹³³ *Id.* at 2-3; *see* Oard et al., *supra* note 9, at 20.

¹³⁴ Hedin et al., *supra* note 9, at 3-4.

¹³⁵ *See id.* at 12-14.

¹³⁶ *Id.* at 8.

¹³⁷ *Id.* at 3.

¹³⁸ *Id.*

- The TA’s opinion, for any document that participants did appeal.

The TREC coordinators used statistical inference to estimate recall, precision, and F_1 for the results each participating team produced.¹³⁹

[36] Assuming participants diligently appealed the first-pass assessments with which they disagreed, it is reasonable to conclude that TREC’s two-pass assessment process yields a reasonably accurate gold standard. Moreover, that same gold standard is suitable to evaluate not only the participants’ submissions, but also the first-pass assessments of the human reviewers.¹⁴⁰

[37] Parts III.A and III.B briefly describe the processes employed by the two participants whose results this Article compares to manual review. Notably, the methods the two participants used differ substantially from those typically described in the industry as “clustering” or “concept search.”¹⁴¹

A. H5 Participation

[38] At TREC 2009, H5 completed one topic (Topic 204).¹⁴² According to Dan Brassil of H5, the H5 process involves three steps: (i) “definition of relevance,” (ii) “partly-automated design of deterministic queries,” and (iii) “measurement of precision and recall.”¹⁴³ “Once relevance is defined, the two remaining processes of (1) sampling and query design and (2) measurement of precision and recall are conducted

¹³⁹ *Id.* at 3, 11-16.

¹⁴⁰ *See* Hedin et al., *supra* note 9, at 13 (describing the construction of the gold standard).

¹⁴¹ *Sedona Search Commentary*, *supra* note 1, at 202-03.

¹⁴² Hedin et al., *supra* note 9, at 6-7.

¹⁴³ E-mail from Dan Brassil to Maura R. Grossman, *supra* note 105.

iteratively – ‘allowing for query refinement and correction’ – until the clients’ accuracy requirements are met.”¹⁴⁴

[39] H5 describes how its approach differs from other information retrieval methods as follows:

It utilizes an iterative issue-focusing and data-focusing methodology that defines relevancy in detail; most alternative processes provide a reductionist view of relevance (e.g.: a traditional coding manual), or assume that different individuals share a common understanding of relevance.

[H5’s approach] is deterministic: each document is assessed against the relevance criteria and a relevant / not relevant determination is made. . . .

[The approach] is built on precision: whereas many alternative approaches start with a small number [of] keywords intended to be broad so as to capture a lot of relevant data (with the consequence of many false positives), H5’s approach is focused on developing in an automated or semi-automated fashion large numbers of deterministic queries that are very precise: each string may capture just a few documents, but nearly all documents so captured will be relevant; and all the strings together will capture most relevant documents in the collection.¹⁴⁵

In the course of its TREC 2009 effort, H5 sampled and reviewed a total of 20,000 documents.¹⁴⁶ H5 declined to quantify the number of person-hours

¹⁴⁴ *Id.*

¹⁴⁵ *Id.* (citing Dan Brassil et al., *The Centrality of User Modeling to High Recall with High Precision Search*, in 2009 IEEE Int’l Conf. on Systems, Man, and Cybernetics, 91, 91-96.

¹⁴⁶ *Id.*

it expended during the seven to eight week time period between the assignment of the topic and the final submission date.¹⁴⁷

B. Waterloo Participation

[40] The University of Waterloo (“Waterloo”) completed four topics (Topics 201, 202, 203, and 207).¹⁴⁸ Waterloo’s approach consisted of three phases: (i) “interactive search and judging,” (ii) “active learning,” and (iii) recall estimation.¹⁴⁹ The interactive search and judging phase “used essentially the same tools and approach [Waterloo] used in TREC 6.”¹⁵⁰ Waterloo coupled the Wumpus search engine¹⁵¹ to a custom web interface that provided document excerpts and permitted assessments to be coded with a single mouse click.¹⁵² Over the four topics, roughly 12,500 documents were retrieved and reviewed, at an average rate of about 3 documents per minute (about 22 seconds per document; 76 hours in

¹⁴⁷ *Id.*; E-mail from Dan Brassil to Maura R. Grossman (Feb. 16, 2011 15:58 EST) (on file with authors).

¹⁴⁸ Cormack & Mojdeh, *supra* 98, at 2.

¹⁴⁹ *Id.* at 1-3.

¹⁵⁰ *Id.* at 2. *See generally*, Gordon V. Cormack et al., *Efficient Construction of Large Test Collections*, in SIGIR '98 PROCEEDINGS OF THE 21ST ANNUAL INTERNATIONAL ACM SIGIR CONFERENCE ON RESEARCH AND DEVELOPMENT IN INFORMATION RETRIEVAL 282, 284 (1998).

¹⁵¹ *Welcome to the Wumpus Search Engine!*, WUMPUS, <http://www.wumpussearch.org/> (last visited Apr. 11, 2011).

¹⁵² *See* Cormack & Mojdeh, *supra* note 98, at 3 & fig.2; *see also infra* Figure 1. “We used the Wumpus search engine and a custom html interface that showed hits-in-context and radio buttons for adjudication Available for reference were links to the full text of the document and to the full email message containing the document, including attachments in their native format.” Cormack & Mojdeh, *supra* note 98, at 3.

total).¹⁵³ Waterloo used the resulting assessments to train an on-line active learning system, previously developed for spam filtering.¹⁵⁴

[41] The active learning system “yields an estimate of the [probability] that each document is relevant.”¹⁵⁵ Waterloo developed an “efficient user interface to review documents selected by this relevance score” (see Figure 2).¹⁵⁶ “The primary approach was to examine unjudged documents in decreasing order of score, skipping previously adjudicated documents.”¹⁵⁷ The process displayed each document as text and, using a single keystroke, coded each document as relevant or not relevant.¹⁵⁸ Among the four topics, “[a]bout 50,000 documents were reviewed, at an average rate of 20 documents per minute (3 seconds per document)” or 42 hours in total.¹⁵⁹ “From time to time, [Waterloo] revisited the interactive search and judging system, to augment or correct the relevance assessments as new information came to light.”¹⁶⁰

¹⁵³ E-mail from Gordon V. Cormack to K. Krasnow Waterman (Feb. 24, 2010 08:25 EST) (on file with authors) (indicating that 12,508 documents were reviewed at a rate of 22 seconds per document, *i.e.*, 76.44 hours in total).

¹⁵⁴ Cormack & Mojdeh, *supra* note 98, at 3.

¹⁵⁵ *Id.* at 3.

¹⁵⁶ *Id.*

¹⁵⁷ *Id.*

¹⁵⁸ *Id.*

¹⁵⁹ Cormack & Mojdeh, *supra* note 98, at 3.

¹⁶⁰ *Id.*

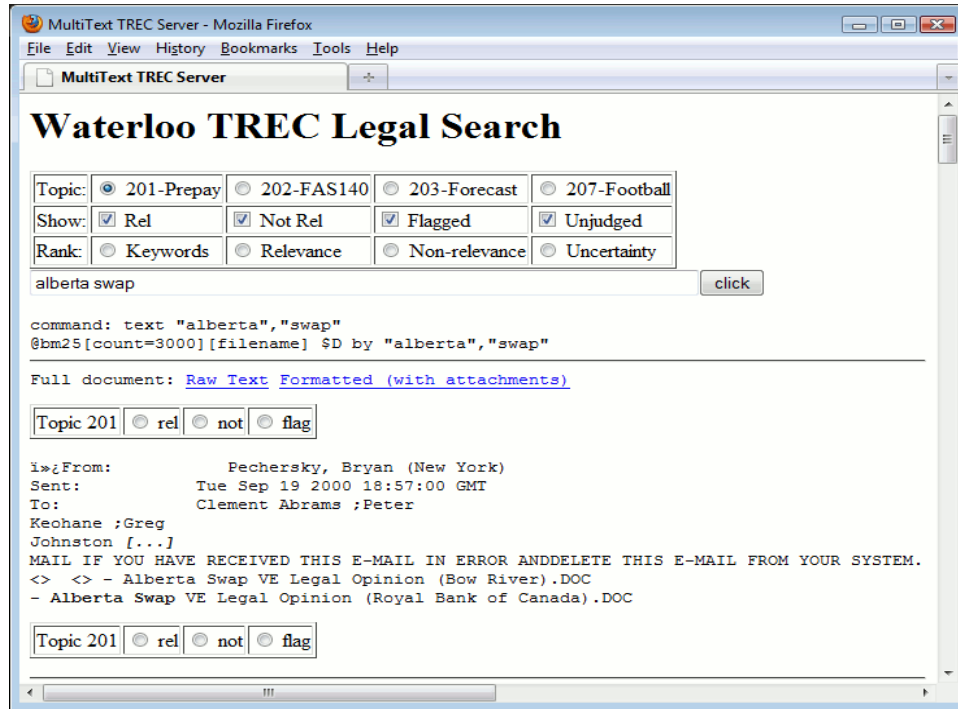


Figure 1: Waterloo's interactive search and judging interface.¹⁶¹

[42] The third and final phase estimated the density of relevant documents as a function of the score assigned by the active learning system, based on the assessments rendered during the active learning phase.¹⁶² Waterloo used this estimate to gauge the tradeoff between recall and precision, and to determine the number of documents to produce so as to optimize F_1 , as required by the task guidelines.¹⁶³

¹⁶¹ *Id.* at 3 & fig.2.

¹⁶² *See id.* at 6.

¹⁶³ *Id.* at 3, 6; *see* Hedin et al., *supra* note 9, at 3.

[43] For Waterloo's TREC 2009 effort, the end result was that a human reviewed every document produced;¹⁶⁴ however, the number of documents reviewed was a small fraction of the entire document population (14,396 of the 836,165 documents were reviewed, on average, per topic).¹⁶⁵ Total review time for all phases was about 118 hours; 30 hours per topic, on average.¹⁶⁶

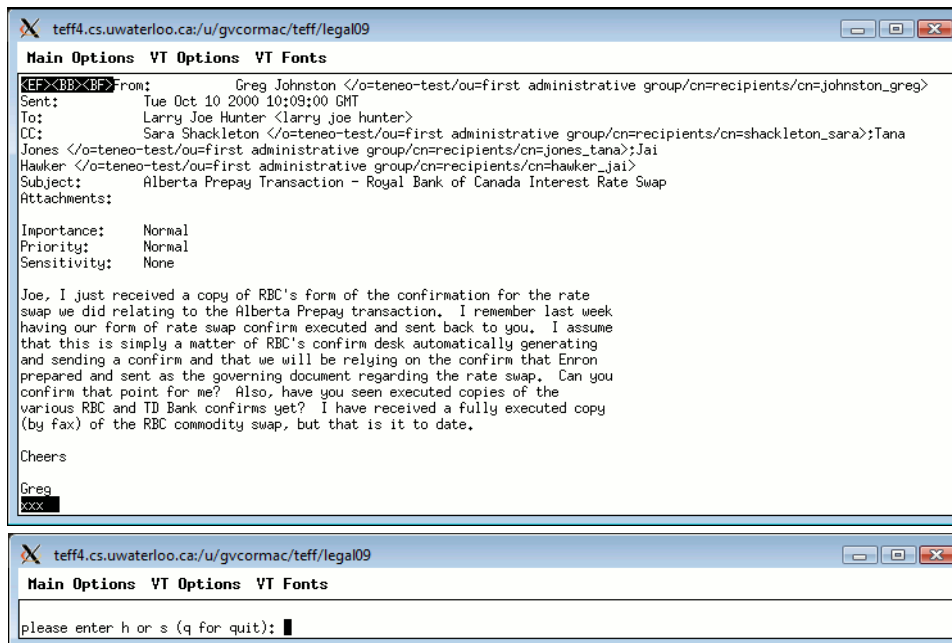


Figure 2: Waterloo's minimalist review interface.¹⁶⁷

¹⁶⁴ See Cormack & Mojdeh *supra* note 98, at 6 (“the optimal strategy was to include *no* unassessed documents”).

¹⁶⁵ *Id.*, at 6 tbl.2; E-mail from Bruce Hedin to Gordon V. Cormack, *supra* note 101 (“I have attached full list of the 836,165 document-level IDs”).

¹⁶⁶ 118 hours is the sum of 76 hours for the interactive search and judging phase (*supra* para. 39) and 42 hours for the active learning phase (*supra* para. 41). Since Waterloo did four topics, the average effort per topic was 29.5 hours.

¹⁶⁷ Cormack & Mojdeh, *supra* note 98, at 4 fig.3.

IV. QUANTITATIVE ANALYSIS

[44] This Article's purpose is to refute the hypothesis that manual review is the best approach by showing that technology-assisted review can yield results that are more nearly complete and more accurate than exhaustive manual review, as measured by recall, precision, and F_1 . To compare technology-assisted to manual review, the study required:

1. The results of one or more technology-assisted reviews. For this purpose, the Authors used the H5 review and the four Waterloo reviews conducted during the course of their participation in the TREC 2009 Legal Track Interactive Task.¹⁶⁸
2. The results of manual reviews for the same topics and datasets as the technology-assisted reviews. For this purpose, the Authors used the manual reviews that TREC conducted on document samples for the purpose of evaluating the results that the participating teams submitted.¹⁶⁹
3. A gold standard determination of relevance or nonrelevance. For this purpose, the Authors used the TREC final adjudicated assessments, for which the TA was the ultimate arbiter.¹⁷⁰

[45] The Authors evaluated the results of the technology-assisted reviews and the manual reviews in exactly the same manner, using the

¹⁶⁸ The TREC results are available online, but use, dissemination and publication of the material is limited. Text REtrieval Conference (TREC), *Past Results*, NAT'L INST. STANDARDS & TECH., <http://trec.nist.gov/results.html> (last visited Apr. 11, 2011) ("Individuals may request access to the protected area containing the raw results by contacting the TREC Program Manager. Before receiving access, individuals will be asked to sign an agreement that acknowledges the limited uses for which the data can be used.").

¹⁶⁹ Text REtrieval Conference (TREC), *Relevance Judgments and Evaluation Tools for the Interactive Task*, NAT'L INST. STANDARDS & TECH., <http://trec.nist.gov/data/legal/09/evalInt09.zip> (last visited Apr. 11, 2011).

¹⁷⁰ *Id.*; see Hedin et al., *supra* note 9, at 2-3.

TREC methodology and the TREC gold standard.¹⁷¹ To compare the effectiveness of the reviews, this Article reports, for each topic:

1. Recall, precision, and F_1 for both the technology-assisted and manual reviews.¹⁷²
2. The *difference* in recall, the difference in precision, and the difference in F_1 between the technology-assisted and manual reviews.¹⁷³
3. The *significance of the difference* for each measure, expressed as P .¹⁷⁴ Traditionally, $P < 0.05$ is interpreted to mean that the difference is statistically significant; $P > 0.1$ is interpreted to mean that the measured difference is not statistically significant. Smaller values of P imply stronger significance; $P < 0.001$ indicates overwhelming significance.¹⁷⁵ The Authors used 100 bootstrap samples of paired differences to estimate the standard error of measurement, assuming a two-tailed normal distribution, to compute P .¹⁷⁶

Table 7 shows recall, precision, and F_1 for the technology-assisted and manual reviews for each of the five topics, as well as the overall average for the five technology-assisted reviews and the five manual reviews. For brevity, the difference in each measure is not shown, but is easily

¹⁷¹ See Hedin et al., *supra* note 9, at 2-5.

¹⁷² See *id.* at 3 (reporting recall, precision, and F_1 for TREC participants); *infra* Table 7 (reporting recall, precision, and F_1 for the TREC manual reviews).

¹⁷³ See *infra* Table 7. A positive difference in some measure indicates that the technology-assisted review is superior in that measure, while a negative difference indicates that it is inferior.

¹⁷⁴ BÜTTCHER ET AL., *supra* note 19, at 426.

¹⁷⁵ See *id.*

¹⁷⁶ See *id.* at 412-31. “The *bootstrap* . . . is a method for simulating an empirical distribution modeling $f(S)$ by sampling the sample s .” *Id.* at 424.

computed from the table. For example, for Topic 201, the difference in recall between Waterloo and TREC is $77.8\% - 75.6\% = +2.2\%$.

| Topic | Team | Recall | Precision | F_1 |
|-------|----------------------|-----------|-----------|-----------|
| 201 | Waterloo | (†) 77.8% | (*) 91.2% | (*) 84.0% |
| | TREC (Law Students) | 75.6% | 5.0% | 9.5% |
| 202 | Waterloo | 67.3% | (*) 88.4% | (*) 76.4% |
| | TREC (Law Students) | (†) 79.9% | 26.7% | 40.0% |
| 203 | Waterloo | (*) 86.5% | (*) 69.2% | (*) 76.9% |
| | TREC (Professionals) | 25.2% | 12.5% | 16.7% |
| 204 | H5 | (*) 76.2% | (*) 84.4% | (*) 80.1% |
| | TREC (Professionals) | 36.9% | 25.5% | 30.2% |
| 207 | Waterloo | 76.1% | (†) 90.7% | 82.8% |
| | TREC (Professionals) | (†) 79.0% | 89.0% | (†) 83.7% |
| Avg. | H5/Waterloo | (†) 76.7% | (*) 84.7% | (*) 80.0% |
| | TREC | 59.3% | 31.7% | 36.0% |

Table 7: Effectiveness of TREC 2009 Legal Track technology-assisted approaches (H5 and Waterloo) compared to exhaustive manual reviews (TREC). Results marked (*) are superior and overwhelmingly significant ($P < 0.001$). Results marked (†) are superior but not statistically significant ($P > 0.1$).¹⁷⁷

[46] For each topic and each measure, the larger value is marked with either (*) or (†); (*) indicates that the measured difference is overwhelmingly significant ($P < 0.001$), while (†) indicates that it is not statistically significant ($P > 0.1$). As Table 7 illustrates, all of the measured differences are either overwhelmingly significant or not statistically significant.¹⁷⁸

V. QUALITATIVE ANALYSIS

[47] The quantitative results show that the recall of the manual reviews varies from about 25% (Topic 203) to about 80% (Topic 202). That is, human assessors missed between 20% and 75% of all relevant documents.¹⁷⁹ Is this shortfall the result of clerical error, a

¹⁷⁷ For the information contained in this table, see *Past Results*, *supra* note 168; *Relevance Judgments and Evaluation Tools for the Interactive Task*, *supra* note 169. For details on the calculation and meaning of P , see *supra* para. 43.

¹⁷⁸ *Supra* Table 7.

¹⁷⁹ *See supra* Table 7.

misinterpretation of relevance, or disagreement over marginal documents whose responsiveness is debatable? If the missed documents are marginal, the shortfall may be of little consequence; but if the missed documents are clearly responsive, production may be inadequate, and under Rule 37(a)(4), such a production could constitute a failure to respond.¹⁸⁰

[48] To address this question, the Authors examined the documents that the TREC assessors coded as nonresponsive to Topics 204 and 207, but H5 or Waterloo coded as responsive, and the TA adjudicated as responsive. Recall from Table 5 that Topic 204 concerned shredding and destruction of documents, while Topic 207 concerned football and gambling. The Authors chose these topics because they were more likely to be easily accessible to the reader, as opposed to other topics, which were more technical in nature. In addition, lawyers employed by professional review companies assessed these two topics using accepted practices for manual review.¹⁸¹

[49] For Topic 204, 160 of the assessed documents were coded as nonresponsive by the manual reviewers and responsive by H5 and the TA;¹⁸² Topic 207, 51 documents met these same criteria except that Waterloo and the TA made the responsiveness determinations.¹⁸³ From these numbers, the Authors extrapolated that the manual reviewers would

¹⁸⁰ See FED. R. CIV. P. 37(a)(4).

¹⁸¹ See Hedin et al., *supra* note 9, at 8 (“The review of the samples for three of the seven Interactive topics (203, 204, and 207) was carried out by two firms that include professional document-review services among their offerings.”).

¹⁸² The Authors identified these documents by comparing the submitted results, *see Past Results*, *supra* note 168 (file input.H52009.gz), the first-pass assessments, *see Relevance Judgments and Evaluation Tools for the Interactive Task*, *supra* note 169 (file qrels_doc_pre_all.txt), and the final adjudicated results, *see id.* (file qrels_doc_post_all.txt).

¹⁸³ The Authors identified these documents by comparing the submitted results, *see Past Results*, *supra* note 168 (file input.watlint.gz), the first-pass assessments, *see Relevance Judgments and Evaluation Tools for the Interactive Task*, *supra* note 169 (file qrels_doc_pre_all.txt), and the final adjudicated results, *see id.* (file qrels_doc_post_all.txt).

have missed 1,918 and 1,273 responsive documents (for Topics 204 and 207, respectively), had they reviewed the entire document collection.

[50] For each of these documents, the Authors used their judgment to assess whether the document had been miscoded due to:

- *Inarguable error*: Under any reasonable interpretation of relevance, the reviewer should have coded the document as responsive, but did not. Possible reasons for such error include fatigue or inattention, overlooking part of the document, poor comprehension, or data entry mistakes in coding the document.¹⁸⁴ For example, a document about “shredding” (see Figure 3) is responsive on its face to Topic 204; similarly “Fantasy Football” (see Figure 4) is responsive on its face to Topic 207.

Date: Tuesday, January 22, 2002 11:31:39 GMT
Subject:

I'm in. I'll be shredding 'till 11am so I should have plenty of time to make it.

Figure 3: Topic 204 Inarguable error. A professional reviewer coded this document as nonresponsive, although it clearly pertains to document shredding, as specified in the production request.¹⁸⁵

¹⁸⁴ Cf. Jeremy M. Wolfe et al., *Low Target Prevalence Is a Stubborn Source of Errors in Visual Search Tasks*, 136 J. EXPERIMENTAL PSYCH. 623, 623-24 (2007) (showing that in visual search tasks, humans have much higher error rates when the prevalence of target items is low).

¹⁸⁵ See *supra* Table 5. Figure 3 is an excerpt from document 0.7.47.1449689 in the TREC 2009 dataset, *supra* note 101.

From: Bass, Eric
Sent: Thursday, January 17, 2002 11:19 AM
To: Lenhart, Matthew
Subject: FFL Dues

You owe \$80 for fantasy football. When can you pay?

Figure 4: Topic 207 Inarguable error. A professional reviewer coded this document as nonresponsive, although it clearly pertains to fantasy football, as specified in the production request.¹⁸⁶

- *Interpretive error*: Under some reasonable interpretation of relevance – but not the TA’s interpretation as provided in the topic guidelines – an assessor might consider the document as nonresponsive. For example, a reviewer might have construed an automated message stating, “your mailbox is nearly full; please delete unwanted messages” (see Figure 5) as nonresponsive to Topic 204, although the TA defined it as responsive. Similarly, an assessor might have construed a message concerning children’s football (see Figure 6) as nonresponsive to Topic 207, although the TA defined it as responsive.

¹⁸⁶ See *supra* Table 5. Figure 4 is an excerpt from document 0.7.47.320807 from the TREC 2009 dataset, *supra* note 101.

WARNING: Your mailbox is approaching the size limit

This warning is sent automatically to inform you that your mailbox is approaching the maximum size limit. Your mailbox size is currently 79094 KB.

Mailbox size limits:

When your mailbox reaches 75000 KB you will receive this message. To check the size of your mailbox:

Right-click the mailbox (Outlook Today),
Select Properties and click the Folder Size button.
This method can be used on individual folders as well.

To make more space available, delete any items that are no longer needed such as Sent Items and Journal entries.

Figure 5: Topic 204 Interpretive error. A professional reviewer coded this automated message as nonresponsive, although the TA construed such messages to be responsive to Topic 204.¹⁸⁷

Subject: RE: Meet w/ Belden

I need to leave at 3:30 today to go to my stepson's football game. Unfortunately, I have a 2:00 and 3:00 meeting already. Is this just a general catch-up discussion?

Figure 6: Topic 207 Interpretive error. The reviewer may have construed a children's league football game to be outside of the scope of "gambling on football." The TA deemed otherwise.¹⁸⁸

- *Arguable error*: Reasonable, informed assessors might disagree or find it difficult to determine whether or not the document met the TA's conception of responsiveness (e.g., Figures 7 and 8).

¹⁸⁷ See *supra* Table 5. Figure 5 is an excerpt from document 0.7.47.1048852 in the TREC 2009 dataset, *supra* note 101.

¹⁸⁸ See *supra* Table 5. Figure 6 is an excerpt from document 0.7.47.668065 in the TREC 2009 dataset, *supra* note 101.

Subject: Original Guarantees

Just a followup note:

We are still unclear as to whether we should continue to send original incoming and outgoing guarantees to Global Contracts (which is what we have been doing for about 4 years, since the Corp. Secretary kicked us out of using their vault on 48 for originals because we had too many documents). I think it would be good practice if Legal and Credit sent the originals to the same place, so we will be able to find them when we want them. So my question to y'all is, do you think we should send them to Global Contracts, to you, or directly to the 48th floor vault (if they let us!).

Figure 7: Topic 204 Arguable error. This message concerns *where* to store particular documents, not specifically their destruction or retention. Applying the TA's conception of relevance, reasonable, informed assessors might disagree as to its responsiveness.¹⁸⁹

Subject: RE: How good is Temptation Island 2

They have some cute guy lawyers this year-but I bet you probably watch that manly Monday night Football.

Figure 8: Topic 207 Arguable error. This message mentions football, but not a specific football team, player, or game. Reasonable, informed reviewers might disagree about whether or not it is responsive according to the TA's conception of relevance.¹⁹⁰

[51] When rendering assessments for the qualitative analysis, the Authors considered the mock complaint,¹⁹¹ the topics,¹⁹² and the topic-specific assessment guidelines memorializing the TA's conception of relevance, which were given to the human reviewers for reference

¹⁸⁹ See *supra* Table 5. Figure 7 is an excerpt from document 0.7.47.1304583 in the TREC 2009 dataset, *supra* note 101.

¹⁹⁰ See *supra* Table 5. Figure 8 shows an excerpt from document 0.7.6.179483 in the TREC 2009 dataset, *supra* note 101.

¹⁹¹ See generally *Complaint, Grumby v. Volteron Corp.*, *supra* note 97.

¹⁹² *Id.* at 14; Hedin et al, *supra* note 9, at 5-6.

purposes.¹⁹³ Table 8 summarizes the findings: The vast majority of missed documents are attributable either to inarguable error or to misinterpretation of the definition of relevance (interpretive error). Remarkably, the findings identify only 4% of all errors as arguable.

| Topic | Error Type | | | Total |
|----------|------------|--------------|----------|-------|
| | Inarguable | Interpretive | Arguable | |
| 204 | 98 | 56 | 6 | 160 |
| 207 | 39 | 11 | 1 | 51 |
| Total | 137 | 67 | 7 | 211 |
| Fraction | 65% | 31% | 4% | 100% |

Table 8: Number of responsive documents that human reviewers missed, categorized by the nature of the error. 65% of missed documents are relevant on their face. 31% of missed documents are clearly relevant, when the topic-specific guidelines are considered. Only 4% of missed documents, in the opinion of the Authors, have debatable responsiveness, according to the topic-specific guidelines.¹⁹⁴

VI. RESULTS AND DISCUSSION

[52] Tables 6 and 7 show that, by all measures, the average efficiency and effectiveness of the five technology-assisted reviews surpasses that of the five manual reviews. The technology-assisted reviews require, on average, human review of only 1.9% of the documents, a fifty-fold savings over exhaustive manual review. For F_1 and precision, the measured difference is overwhelmingly statistically significant ($P < 0.001$);¹⁹⁵ for recall the measured difference is not significant ($P > 0.1$).¹⁹⁶ These measurements provide strong evidence that the technology-assisted

¹⁹³ Text REtrieval Conference (TREC), *TREC-2009 Legal Track – Interactive Task, Topic-Specific Guidelines – Topic 204*, U. WATERLOO, http://plg1.cs.uwaterloo.ca/trec-assess/TopicGuidelines_204.pdf (last updated Oct. 22, 2009); Text REtrieval Conference (TREC), *TREC-2009 Legal Track – Interactive Task, Topic-Specific Guidelines – Topic 207*, U. WATERLOO, http://plg1.cs.uwaterloo.ca/trec-assess/TopicGuidelines_207_.pdf (last updated Oct. 22, 2009).

¹⁹⁴ See sources cited *supra* note 193.

¹⁹⁵ See *supra* Tables 6, 7.

¹⁹⁶ *Id.*

processes studied here yield better overall results, and better precision, in particular, than the TREC manual review process. The measurements also suggest that the technology-assisted processes may yield better recall, but the statistical evidence is insufficiently strong to support a firm conclusion to this effect.

[53] It should be noted that the objective of TREC participants was to maximize F_1 , not recall or precision, per se.¹⁹⁷ It happens that they achieved, on average, higher precision.¹⁹⁸ Had the participants considered recall to be more important, they might have traded off precision (and possibly F_1) for recall, by using a broader interpretation of relevance, or by adjusting a sensitivity parameter in their software.

[54] Table 7 shows that, for four of the five topics, the technology-assisted processes achieve substantially higher F_1 scores, largely due to their high precision. Nonetheless, for a majority of the topics, the technology-assisted processes achieve higher recall as well; for two topics, substantially higher.¹⁹⁹ For Topic 207, there is no meaningful difference in effectiveness between the technology-assisted and manual reviews, for any of the three measures. *There is not one single measure for which manual review is significantly better than technology-assisted review.*

[55] For three of the five topics (Topics 201, 202, and 207) the results show no significant difference in recall between the technology-assisted and manual reviews. This result is perhaps not surprising, since the recall scores are all on the order of 70% – the best that might be reasonably achieved, given the level of agreement among human assessors. As such, the results support the conclusion that technology-assisted review can achieve at least as high recall as manual review, and higher precision, at a fraction of the review effort, and hence, a fraction of the cost.

¹⁹⁷ See Hedin et al., *supra* note 9, at 15.

¹⁹⁸ See *supra* Tables 6, 7.

¹⁹⁹ See *supra* Table 7.

VII. LIMITATIONS

[56] The 2009 TREC effort used a mock complaint and production requests composed by lawyers to be as realistic as possible.²⁰⁰ Furthermore, the role of the TA was intended to simulate that of a senior attorney overseeing a real document review.²⁰¹ Finally, the dataset consisted of real e-mail messages captured within the context of an actual investigation.²⁰² These components of the study are perhaps as realistic as might reasonably be achieved outside of an actual legal setting.²⁰³ One possible limitation is that the Enron story, and the Enron dataset, are both well known, particularly since the Enron documents are frequently used in vendor product demonstrations.²⁰⁴ Both participants and TAs may have had prior knowledge of both the story and dataset, affecting their strategies and assessments. In addition, there is a tremendous body of extrinsic information that may have influenced participants and assessors alike, including the results of the actual proceedings, commentaries,²⁰⁵ books,²⁰⁶

²⁰⁰ Hedin et al., *supra* note 9, at 2.

²⁰¹ *See id.*; *see also* Oard et al., *supra* note 9, at 20.

²⁰² *See* Hedin et al., *supra* note 9, at 4.

²⁰³ *See id.*

²⁰⁴ *See, e.g.*, John Markoff, Armies of Expensive Lawyers Replaced by Cheaper Software, N.Y. TIMES, Mar. 5, 2011, A1, available at <http://www.nytimes.com/2011/03/05/science/05legal.html>; *see also* E-mail from Jonathan Nystrom to Maura R. Grossman (Apr. 5, 2011 19:12 EDT) (on file with authors) (confirming use of the Enron data set for product demonstrations); E-mail from Jim Renehan to Maura R. Grossman (Apr. 5, 2011 20:06 EDT) (on file with authors) (confirming use of the Enron data set for product demonstrations); E-mail from Lisa Schofield to Maura R. Grossman (Apr. 5, 2011 18:27 EDT) (on file with authors) (confirming use of the Enron data set for product demonstrations); E-mail from Edward Stroz to Maura R. Grossman (Apr. 5, 2011 18:32 EDT) (on file with authors) (confirming use of the Enron data set for product demonstrations).

²⁰⁵ *See, e.g.*, John C. Coffee Jr., *What Caused Enron?: A Capsule Social and Economic History of the 1990's*, 89 CORNELL L. REV. 269 (2004); Paul M. Healy & Krishna G. Palepu, *The Fall of Enron*, 17 J. ECON. PERSP. 3 (2003).

and even a popular movie.²⁰⁷ It is unclear what effect, if any, these factors may have had on the results.

[57] In general, the TREC teams were privy to less detailed guidance than the manual reviewers, placing the technology-assisted processes at a disadvantage. For example, Topic 202 required the production of documents related to “transactions that the Company characterized as compliant with FAS 140.”²⁰⁸ Participating teams were required to undertake research to identify the relevant transactions, as well as the names of the parties, counterparties, and entities involved.²⁰⁹ Manual reviewers, on the other hand, were given detailed guidelines specifying these elements.²¹⁰

[58] Moreover, TREC conducted manual review on a stratified sample containing a higher proportion of relevant documents than the collection as a whole,²¹¹ and used statistical inference to evaluate the result of reviewing every document in the collection.²¹² Beyond the statistical uncertainty, there also is uncertainty as to whether manual reviewers would have had the same error rate had they reviewed the entire collection. It is not unreasonable to think that, because the proportion of relevant documents would have been lower in the collection than it was in the sample, reviewer recall and precision might have been even lower, because reviewers would have tended to miss the needles in the haystacks due to fatigue, inattention, boredom, and related human factors. This

²⁰⁶ See, e.g., LOREN FOX, ENRON: THE RISE AND FALL (2002); BETHANY MCLEAN AND PETER ELKIND, THE SMARTEST GUYS IN THE ROOM: THE AMAZING RISE AND SCANDALOUS FALL OF ENRON (2003).

²⁰⁷ ENRON: THE SMARTEST GUYS IN THE ROOM (Magnolia Pictures 2005).

²⁰⁸ Hedin et al., *supra* note 9, at 5.

²⁰⁹ See *id.* at 8.

²¹⁰ See *id.* at 3.

²¹¹ See *id.* at 12, tbl.3.

²¹² See generally *id.*

sampling effect, combined with the greater guidance provided to the human reviewers, may have resulted in an overestimate of the effectiveness of manual review, and thus understated the results of the study.

[59] Of note is the fact that the appeals process involved reconsideration – and potential reversal – *only* of manual coding decisions that one or more participating teams appealed, presumably because their results disagreed with the manual reviewers’ decisions.²¹³ The appeals process depended on participants exercising due diligence in identifying the assessments with which they disagreed.²¹⁴ And while it appears that H5 and Waterloo exercised such diligence, it became apparent to the Authors during the course of their analysis that a few assessor errors were overlooked.²¹⁵ These erroneous assessments were deemed correct under the gold standard, with the net effect of overstating the effectiveness of manual reviews, while understating the effectiveness of technology-assisted review.²¹⁶ It is also likely that the manual review and technology-assisted processes incorrectly coded some documents that were not appealed.²¹⁷ The impact of the resulting errors on the gold standard would be to overstate both recall and precision for manual review, as well as for technology-assisted review, with no net advantage to either.

²¹³ See Hedin et al., *supra* note 9 at 3, 13-14. There is no benefit, and therefore no incentive, for participating teams to appeal coding decisions with which they agree.

²¹⁴ See *id.* If participating teams do not appeal the manual reviewers’ incorrect decisions, those incorrect decisions will be incorporated into the gold standard, compromising its accuracy and usefulness.

²¹⁵ Hedin et al., *supra* note 9 at 14, tbl.4 (showing that for every topic, H5 and Waterloo appealed the majority of disagreements between their results and the manual assessments).

²¹⁶ See *supra* note 214. If the manual review is incorrect, and the technology-assisted review is correct, the results will overstate the effectiveness of manual review at the expense of technology-assisted review.

²¹⁷ Given that neither the manual reviewers nor the technology-assisted processes are infallible, it stands to reason that they may occasionally agree on coding decisions that are incorrect.

[60] In designing this study, the Authors considered only the results of two of the eleven teams participating in TREC 2009, because they were considered most likely to demonstrate that technology-assisted review can improve upon exhaustive manual review. The study considered all submissions by these two teams, which happened to be the most effective submissions for five of the seven topics. The study did not consider Topics 205 and 206, because neither H5 nor Waterloo submitted results for them. Furthermore, due to a dearth of appeals, there was no reliable gold standard for Topic 206.²¹⁸ The Authors were aware before conducting their analysis that the H5 and Waterloo submissions were the most effective for their respective topics. To show that the results are significant in spite of this prior knowledge, the Authors applied Bonferroni correction,²¹⁹ which multiplies P by 11, the number of participating teams. Even under Bonferroni correction, the results are overwhelmingly significant.

VIII. CONCLUSION

[61] Overall, the myth that exhaustive manual review is the most effective – and therefore, the most defensible – approach to document review is strongly refuted. Technology-assisted review can (and does) yield more accurate results than exhaustive manual review, with much lower effort. Of course, not all technology-assisted reviews (and not all manual reviews) are created equal. The particular processes found to be superior in this study are both interactive, employing a combination of computer and human input. While these processes require the review of orders of magnitude fewer documents than exhaustive manual review, neither entails the naïve application of technology absent human judgment. Future work may address *which* technology-assisted review process(es) will improve *most* on manual review, not *whether* technology-assisted review *can* improve on manual review.

²¹⁸ Hedin et al., *supra* note 9, at 17-18 (“Topic 206 represents the one topic, out of the seven featured in the 2009 exercise, for which we believe the post-adjudication results are not reliable. . . . We do not believe, therefore, that any valid conclusions can be drawn from the scores recorded for this topic . . .”).

²¹⁹ See BÜTTCHER ET AL., *supra* note 19, at 428.