

Coopetition in IR Research

Ellen M. Voorhees

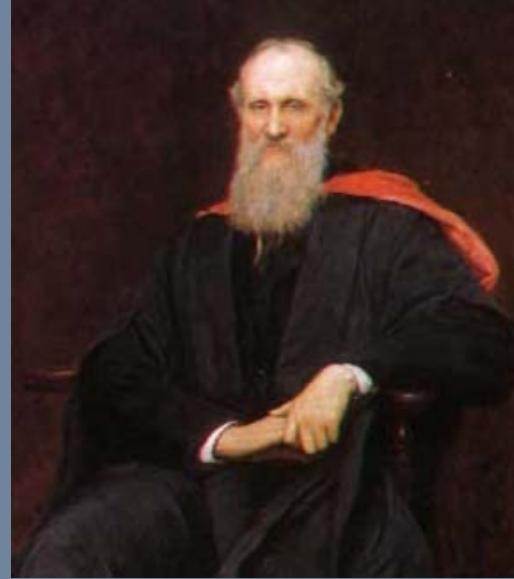


National Institute of Standards and Technology



To promote U.S. innovation and industrial competitiveness by advancing measurement science, standards, and technology in ways that enhance economic security and improve our quality of life.

NIST mission statement



If you can not measure it, you can not improve it.

Lord Kelvin

Dr. Kelvin

Symbol: K



SI unit of
TEMPERATURE

©NIST

Candela

Symbol: cd



SI unit of
BRIGHTNESS

©NIST

Monsieur Kilogram

Symbol: kg



SI unit of
MASS

©NIST

Major Uncertainty

Symbol: ?



ERROR!

©NIST

Meter Man

Symbol: m



SI unit of
LENGTH

©NIST

The Mole

Symbol: mol

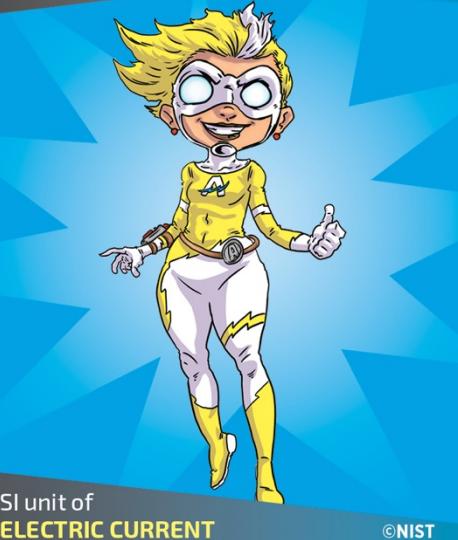


SI unit of
AMOUNT OF SUBSTANCE

©NIST

Ms. Ampere

Symbol: A



SI unit of
ELECTRIC CURRENT

©NIST

Professor Second

Symbol: s



SI unit of
TIME

©NIST

Human Language Technology Measurement

Ida Rhodes and the Problem with 'Water Goats'

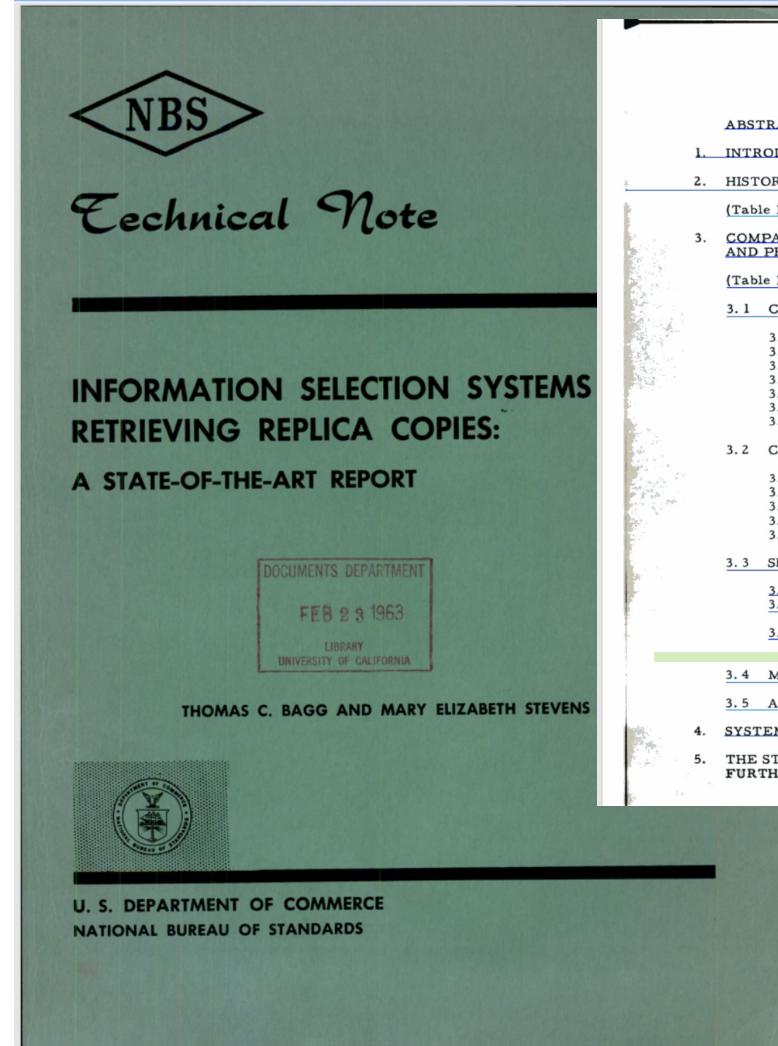
BY KATIE RAPP ON MARCH 16, 2016

ELECTRONICS, INFORMATION TECHNOLOGY, MATHEMATICS AND STATISTICS, NIST GENERAL



Mikhail Dudarev/Fotolia/EpicStockMedia/Shutterstock/Hanacek/NIST

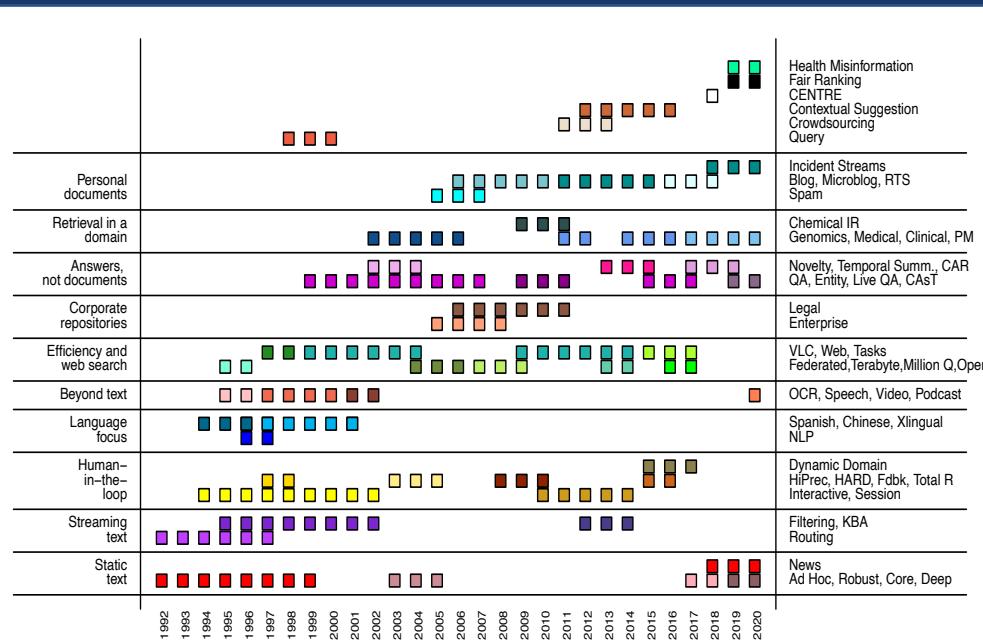
NIST (NBS) has a long history of language evaluation work



Adolf

Text REtrieval Conference (TREC)

Series of community evals (“coopetitions”) that build research infrastructure.



pioneered use of “pooling” for building large collections

built > 150 test collections for dozens of search tasks

hundreds of participant teams world-wide

premier venue for determining research methodology

Coopetition?



Competing may give you a bigger piece of the pie ...

- Benchmark: standardized task
- Coopetition: “cooperative competition” (Wikipedia)
- In IR, built through events such as TREC, CLEF, NTCIR, FIRE, INEX...

Coopetition?



...while cooperation makes
the whole pie bigger.

- Wikipedia: “cooperative competition”
- In IR, known by a variety of other names:
 - challenge problems
 - community evaluations
 - shared tasks
- Exemplified by TREC, CLEF, NTCIR, FIRE, INEX...

Cranfield Paradigm



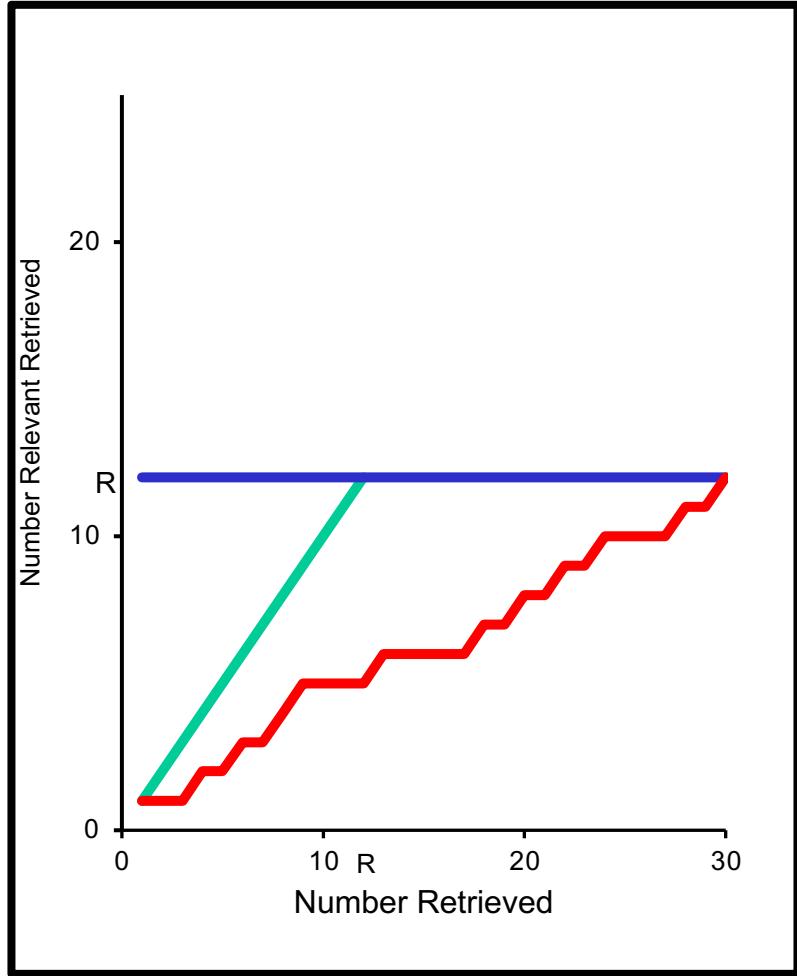
a test collection



Cyril Cleverdon

- Laboratory testing of retrieval systems first done in Cranfield II experiment (1963)
 - fixed document and query sets
 - evaluation based on relevance judgments
- Test collections
 - set of documents
 - set of questions
 - relevance judgments

Cranfield Paradigm



- Retrieval system response to a question is a ranked list of documents.
- The ideal output is a list with all relevant documents ranked before any non-relevant document.
- Easy to compute a variety of different evaluation measures from a ranked list if you know the set of relevant documents

Cranfield Paradigm



Defines “core competency” of IR:
retrieve relevant before non-relevant



Different effectiveness metrics provide different abstractions of real-world task(s).



Abstracted task is a necessary but not sufficient proxy for real task.



Significance tests suggest whether a change is worthwhile.

What is a good search result?

TREC
About 13,900,000 results (0.14 seconds)
Advanced search

- + TREC - Home Page www.trec.state.tx.us/ - Cached
Official site of the Texas Real Estate Commission, the body that governs real estate practices in the state. Includes licensing information, laws, contact ...
 - TREC - Licenses Main Page www.trec.state.tx.us/licenses/default.asp
Licensing Area Main Page ...
 - My License Online Services mylicense.trec.state.tx.us/
Texas Real Estate Commission ...
 - Forms, Laws, Contracts Main Page www.trec.state.tx.us/formslawscontrac...
TREC Forms in Adobe Acrobat ...
 - More results from trec.state.tx.us >
- Ted REtrieval Conference (TREC) Home Page trec.net.gov/ - Cached
Aug 1, 2000 - An annual information retrieval conference and competition, the purpose of which is to support and further research within the information ...
- Tennessee Real Estate Commission - Home www.tn.gov/commerce/boards/trec - Cached
The TREC meetings will be held in Chattanooga in September, Kingsport in October, and ... Please contact the TREC office to preregister to attend the meeting. ...
- The Real Estate Council - Home Page www.reccouncil.com/ - Cached
The Giving Gala, presented by Deloitte, celebrates The Real Estate Council Foundation's good works in four key grant areas: housing, education, job creation ...
- TREC Rentals www.trecrental.com/ - Cached
ROOT! Drive-In 3 Photo Studios in the heart of Chelsea Drive In and Water Shoot Capable.
- TREC - Training Resources for the Environmental Community www.trec.org/ - Cached



Search is inherently a user activity

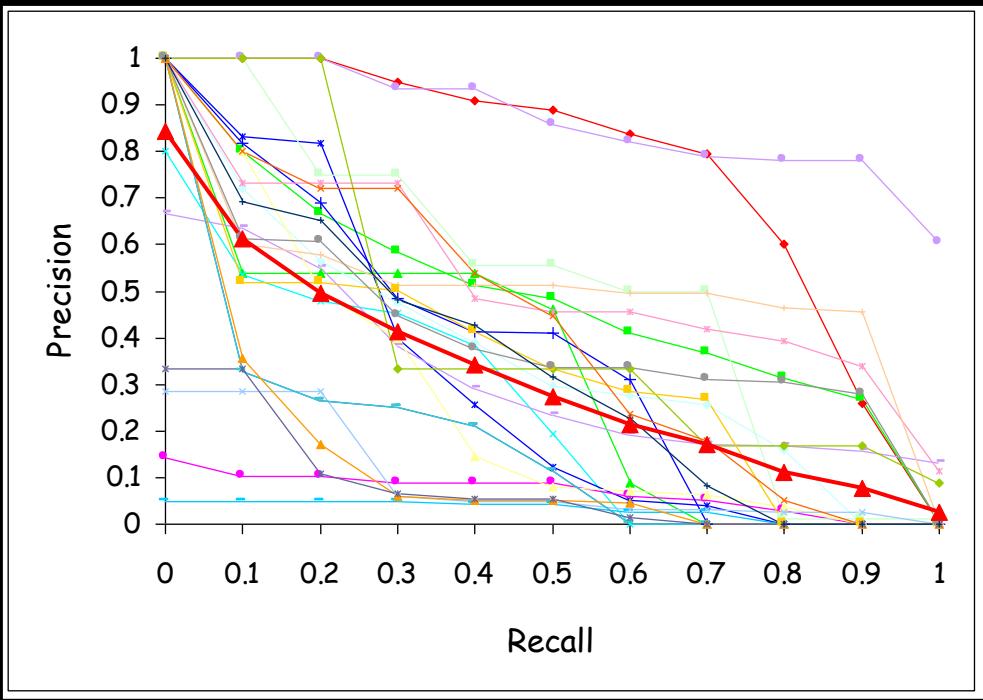


Different users have different (conflicting) criteria for success



There is no single Truth

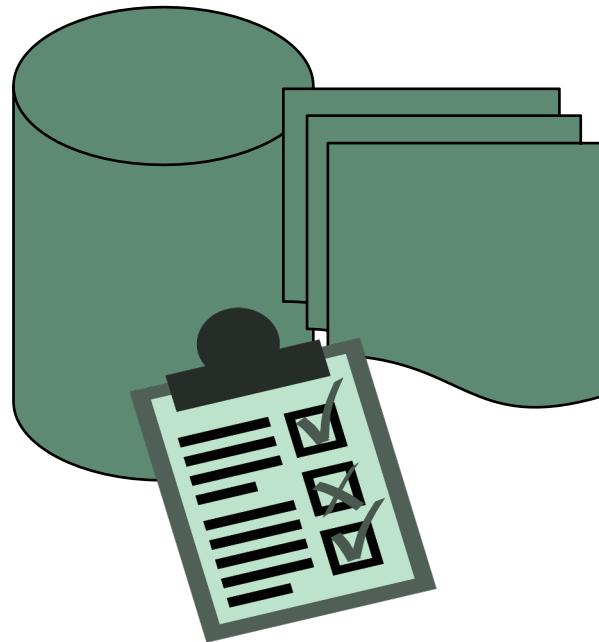
Scoring Ranked Lists



Interpolated precision at standard recall points for 25 individual topics from the best TREC-7 ad hoc run, plus the average curve over all 50 topics in TREC-7 test set for that run (heavy red line).

- Variance across topics is large---larger than differences between systems
- Signal to be measured is small relative to noise
- This variability is why Cranfield averages scores over many topics

Rationale for Cranfield

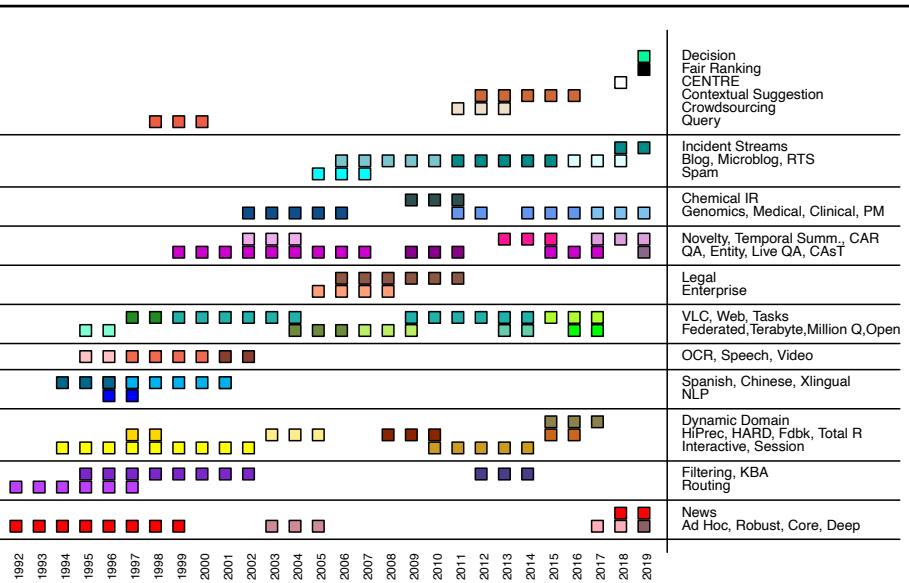


Sufficient fidelity to real user tasks to be informative

General enough to be broadly applicable, feasible, relatively inexpensive

Lose realism to gain control over variables: more experimental power at lower cost

Competitions in IR

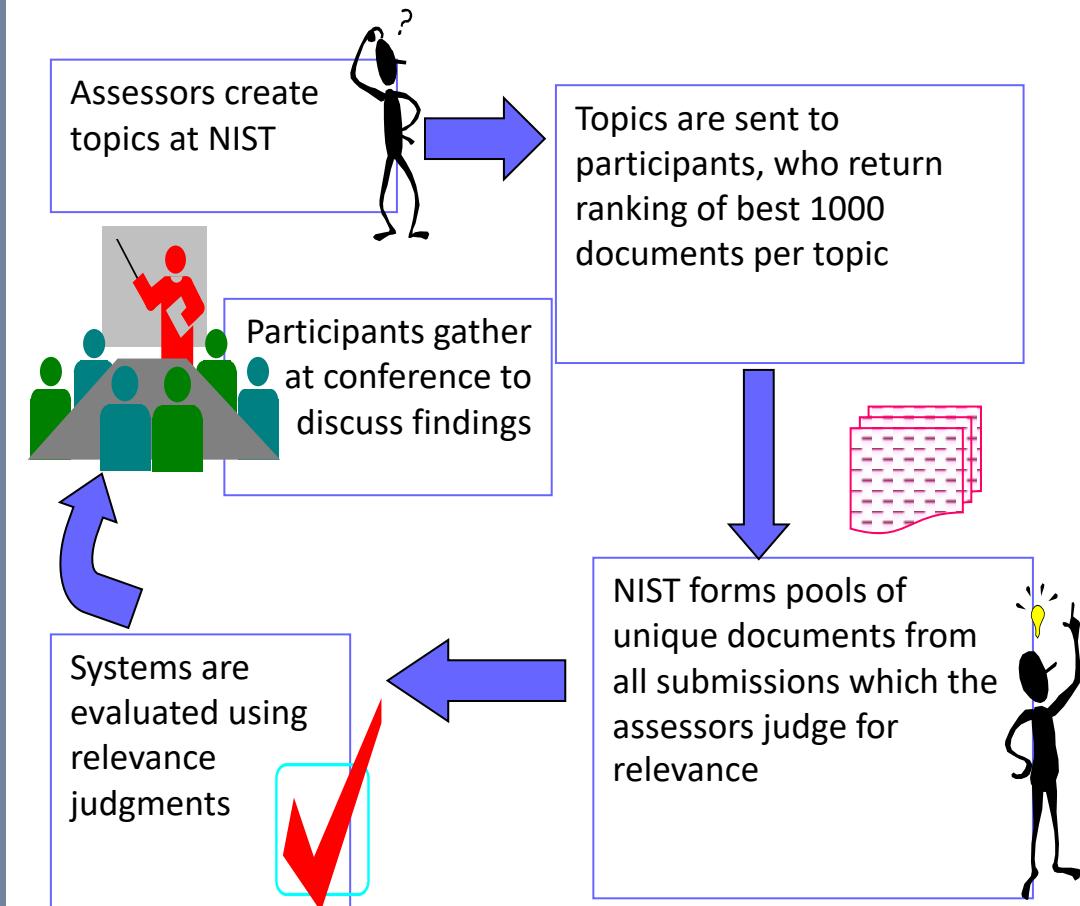


TREC tracks through the years

- Competitions are where the infrastructure that enables Cranfield is created
- Infrastructure includes
 - test collections
 - task definitions (with measures)
 - methodology

Competitions in IR

- Organizers provide appropriate data set, generally “documents” and “topics”.
- Participants use their own systems to create results that are returned to organizers.
- Some of the results get annotated by humans.
- Annotations used to score results; scores returned to participants.
- Data, annotations, results, scores archived for future use.



Benefits of Coopetitions

Improve the state-of-the-art

Establish the research methodology

Form/solidify a research community

Facilitate technology transfer

Amortize the costs of infrastructure

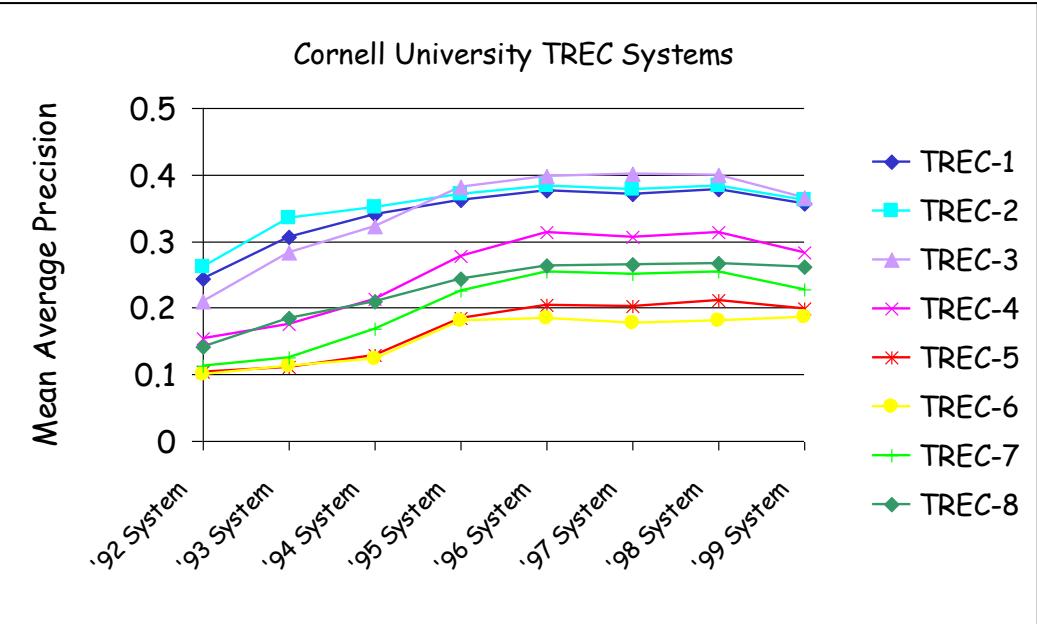
Risks of Coopetitions

Community overfitting
to single dataset/task

Poor task
abstraction

Method conformity

State of the Art



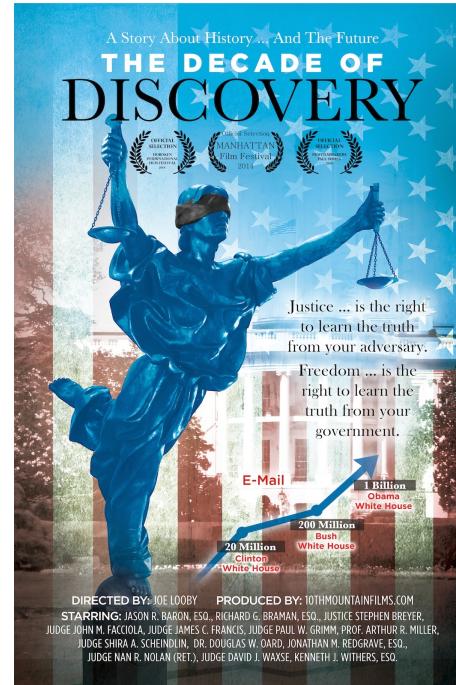
The TREC data revitalized research on information retrieval. Having a standard, widely available, and carefully constructed set of data laid the groundwork for further innovation in the field. The yearly TREC conference fostered collaboration, innovation, and a measured dose of competition (and bragging rights) that led to better information retrieval.

Hal Varian, Google Chief Economist
March 4, 2008

- Improve and document the state of the art
 - scores doubled in early years of TREC
 - improvement driven largely by the test collections constructed
- Risk: entire community overfits
 - construct multiple collections for a task
 - evolve task

Research Community

- Infrastructure enables research in an area, which attracts critical mass
- Risk: poor task abstraction
- Good eval task:
 - abstraction of real-world task so variables affecting performance can be controlled...
 - ...but must capture salient aspects of real task or exercise is pointless
 - metrics must accurately predict relative effectiveness on those aspects
 - adequate level of difficulty
 - best if measures are diagnostic



This project [the TREC Legal track] can be expected to identify both cost effective and reliable search and information retrieval methodologies and best practice recommendations, which, if adhered to, certainly would support an argument that the party employing them performed a reasonable ESI search, whether for privilege review or other purposes.

Magistrate Judge Paul W. Grimm
Victor Stanley v. Creative Pipe

Research Methodology

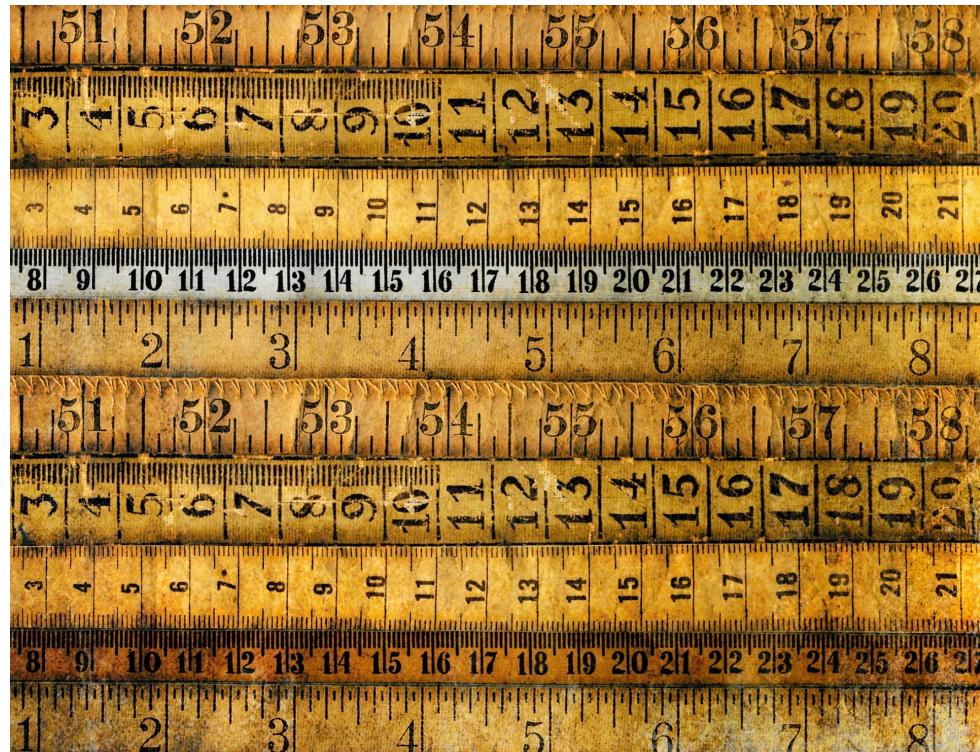


Image: arielrobin/Pixabay

TREC is an annual benchmarking exercise that has become a de facto standard in Information Retrieval evaluation.

- Way to do research in area
 - measures
 - experimental design
- Standardize evaluation measures
 - e.g., interpolation
 - e.g., measure parameter values
- [Risk: data set abuse
 - e.g., compare MAP on TREC web collections
 - e.g., relevance == popularity]

Technology Transfer



TREC has proven to be a valuable forum in which IBM Research has contributed to an improved understanding of search, while at the same time the insights obtained by participating in TREC have helped to improve IBM's products and services.

Alan Marwick, et al.
IBM chapter of the TREC book
2005

- Effective approaches adapted and incorporated into variety of methods
 - e.g., weighting schemes
 - transfer among all sectors
- Risk: conformity of techniques
 - don't want to standardize too early

Costs of Infrastructure

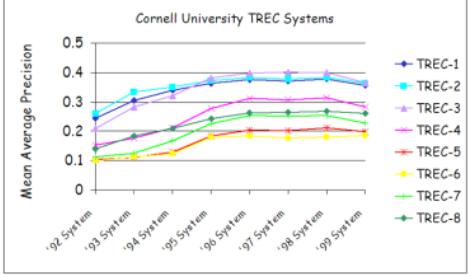
- Building infrastructure does require resources
 - financial
 - researcher time
- Shared infrastructure amortizes the costs over entire community
- [Risk: costs of participation
 - individual can save participation costs assuming others participate so infrastructure gets built
 - minimize competition-only costs such as reporting formats]



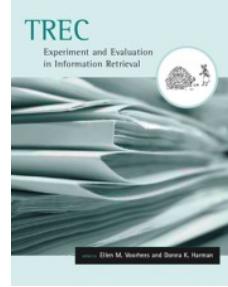
In other words, for every \$1 NIST and its partners invested in TREC, at least \$3.35 to \$5.07 in benefits accrued to IR researchers...These responses [to RTI's survey of IR researchers] suggest that the benefits of TREC to both private and academic organizations go well beyond those quantified by this study's economic benefits.

RTI International
**Economic Impact Assessment
of NIST's Text Retrieval
Conference (TREC) Program**
December 2010

Coopetitions



Improve the state-of-the-art



Establish the research methodology



Form/solidify a research community



Facilitate technology transfer

...for every \$1 invested in TREC, at least \$3.35 to \$5.07 in benefits accrued...



Amortize the cost of infrastructure

A good benchmark task:

- abstraction of real-world task so variables affecting performance can be controlled...
- ...but must capture salient aspects of real task or exercise is pointless
- metrics must accurately predict relative effectiveness on those aspects
- adequate level of difficulty
- best if measures are diagnostic