

[https://commons.wikimedia.org/wiki/File:Longshoreman,\\_Brooklyn,\\_ca.\\_1872-1887.\\_\(5832938139\).jpg](https://commons.wikimedia.org/wiki/File:Longshoreman,_Brooklyn,_ca._1872-1887._(5832938139).jpg)

# Repeatable Runs for Test Collection Documentation

Ian Soboroff, NIST

---

*SIGIR 2019 OSIRRC Workshop, Paris, France*



# Data



[TREC home](#)



---

[Versions of trec\\_eval](#)

[Ad hoc Test Collections](#)

[Web Test Collections](#)

[Blog Track](#)

[Chemical IR Track](#)

[Clinical Decision Support Track](#)

[Common Core Track](#)

[Confusion Track](#)

[Contextual Suggestion Track](#)

[Crowdsourcing Track](#)

[Dynamic Domain Track](#)

[Enterprise Track](#)

[Entity Track](#)

[Filtering Track](#)

# TRECipedia: the Digital Library of IR Test Collections

Browse test collections by...

## Data type

- Newswire
- Web
- Tweets
- Blog
- Medical Cases
- ... other ...

## Document set

- ClueWeb12
- KBA StreamCorpus
- TREC CDs 4 and 5 (minus CR)
- Tweets2011
- ... other ...

## Conference and year

- TREC 2016
- 1999 (which was TREC 8)

## Search task

- **adhoc** search
- filtering / routing
- question answering
- summarization

Popular test collections

- TREC 8 **adhoc**
- TREC 2012 **web**




Most recently updated collections:

- TREC 2004 Robust
- TREC 7 filtering
- trec\_2018\_core
- trec\_2017\_core
- trec\_4\_database\_merging
- trec\_1\_routing
- trec\_1\_adhoc
- trec\_2\_routing
- trec\_2\_adhoc
- trec\_2011\_medical\_records

Links:

-  <http://www.itl.nist.gov>
-  <http://www.nist.gov>
-  <http://www.commerce.gov>

-  <http://trec.nist.gov>
-  <http://ir.nist.gov>

-  [Privacy policy, security notice, accessibility statement](#)
-  [Disclaimer](#)
-  [Freedom of Information Act](#)

## TREC 2004 Robust

### Test Collection Metadata

<b>Conference:</b>	trec_2004
<b>Year:</b>	2004
<b>Document collection:</b>	cd45_minus_cr
<b>Task:</b>	adhoc
<b>Topic set:</b>	robust04
<b>Data type:</b>	newswire
<b>Topics file:</b>	 <a href="https://trec.nist.gov/data/robust/04.testset.gz">https://trec.nist.gov/data/robust/04.testset.gz</a>
<b>Relevance judgments:</b>	 <a href="https://trec.nist.gov/data/robust/qrels.robust2004.txt">https://trec.nist.gov/data/robust/qrels.robust2004.txt</a>
<b>Overview paper:</b>	 <a href="https://trec.nist.gov/pubs/trec13/papers/ROBUST.OVERVIEW.pdf">https://trec.nist.gov/pubs/trec13/papers/ROBUST.OVERVIEW.pdf</a>
<b>Evaluation scores:</b>	 <a href="https://trec.nist.gov/pubs/trec13/appendices/robust.results.html">https://trec.nist.gov/pubs/trec13/appendices/robust.results.html</a>
<b>Top score:</b>	0.333
<b>Median score:</b>	0.2561
<b>Main measure:</b>	map
<b>main_condition:</b>	automatic, title only, 249 topics
<b>Baseline score:</b>	0.2903
<b>Baseline link:</b>	 <a href="https://github.com/osirrc/anserini-docker/tree/v0.1.1">https://github.com/osirrc/anserini-docker/tree/v0.1.1</a>
<b>Baseline notes:</b>	BM25+RM3 model

The robust retrieval track explores methods for improving the consistency of retrieval technology by focusing on poorly performing topics. The retrieval task in the track is a traditional ad hoc retrieval task where the evaluation methodology emphasizes a system's least effective topics. The most promising approach to improving poorly performing topics is exploiting text collections other than the target collection such as the web.

The 2004 edition of the track used 250 topics and required systems to rank the topics by predicted difficulty. The 250 topics within the test set allowed the stability of evaluation measures that emphasize poorly performing topics to be investigated. A new measure, a variant of the traditional MAP measure that uses a geometric mean rather than an arithmetic mean to average individual topic results, shows promise of giving appropriate emphasis to poorly performing topics while being more stable at equal topic set sizes.



**Overview paper:**  <https://trec.nist.gov/pubs/trec13/papers/ROBUST.OVERVIEW.pdf>

**Evaluation scores:**  <https://trec.nist.gov/pubs/trec13/appendices/robust.results.html>

**Top score:** 0.333

**Median score:** 0.2561

**Main measure:** map

**main\_condition:** automatic, title only, 249 topics

**Baseline score:** 0.2903

**Baseline link:**  <https://github.com/osirrc/anserini-docker/tree/v0.1.1>

**Baseline notes:** BM25+RM3 model

The robust retrieval track explores methods for improving the consistency of retrieval techniques. The retrieval task in the track is a traditional ad hoc retrieval task where the evaluation metric is the mean average precision (MAP). The most promising approach to improving poorly performing topics is exploiting text collections.

The 2004 edition of the track used 250 topics and required systems to rank the topics by relevance. The track allowed the stability of evaluation measures that emphasize poorly performing topics to be improved. The MAP measure that uses a geometric mean rather than an arithmetic mean to average individual topic scores gives more emphasis to poorly performing topics while being more stable at equal topic set sizes.



# My questions for OSIRRC...

---

- ❖ Can we do better at linking the image to the source code? What kind of reuse do we want to enable?
- ❖ Can we document what kind of resources are needed?
- ❖ How about pre-indexed collections in a standoff Docker for really big collections like ClueWeb?