



Entity-oriented Filtering of Large Streams

TREC KBA 2013

Knowledge Base Acceleration

John R. Frank

jrf@mit.edu

Ian Soboroff

ian.soboroff@nist.gov

Ellen Voorhees

ellen.voorhees@nist.gov

Max Kleiman-Weiner

maxkw@mit.edu

Dan A. Roberts

drob@mit.edu

Steven J. Bauer

bauer@alum.mit.edu

Nilesh Tripuraneni

nilesht@mit.edu

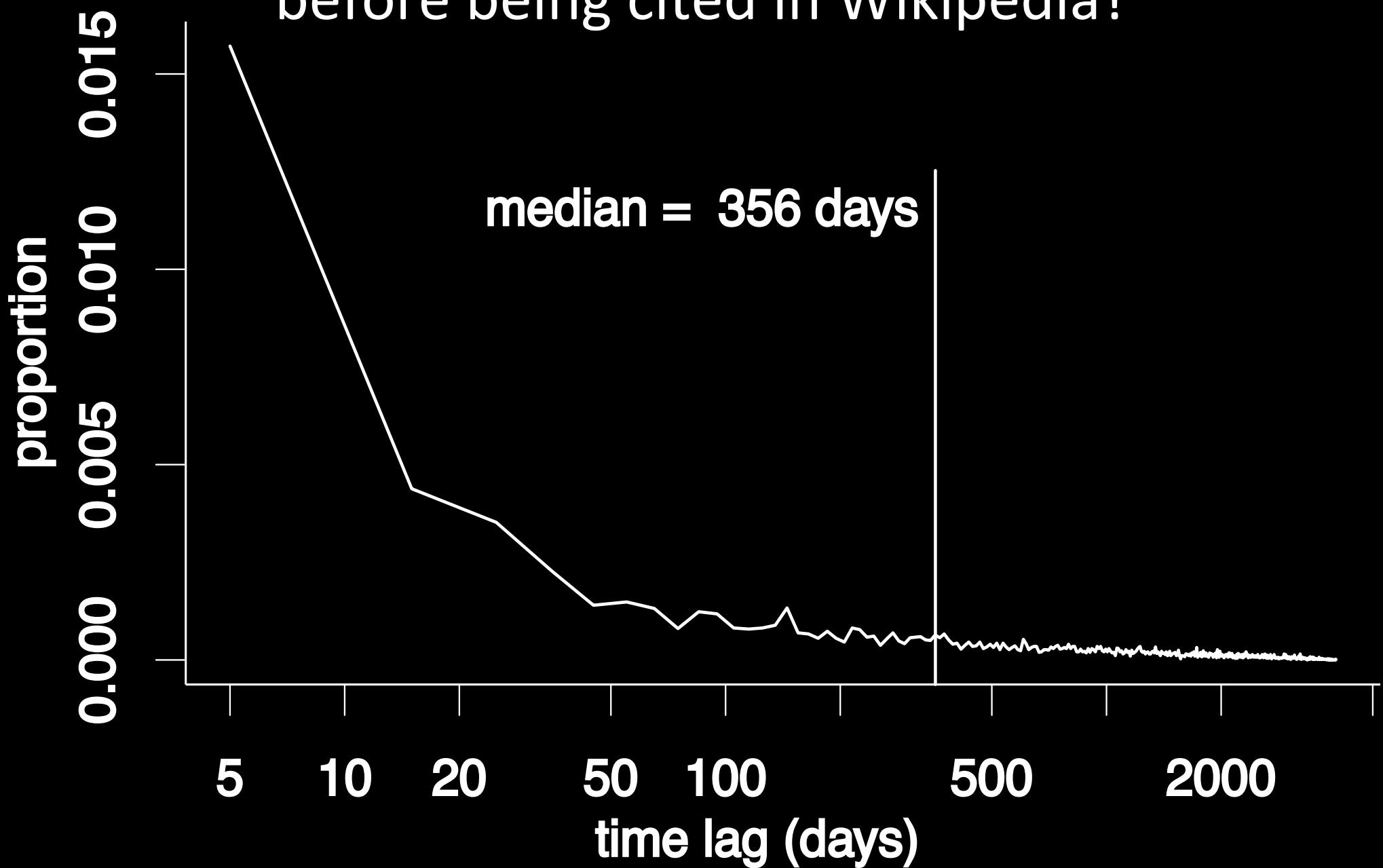
Chris Ré

chrismre@cs.stanford.edu

Ce Zhang

czhang@cs.wisc.edu

How many days must a news article wait before being cited in Wikipedia?



Accelerate?

rate of assimilation << stream size

editors << # entities << # mentions

(definition of a “large” KB)

Date: Tue, 13 Mar 2012 02:45:40 +0000

From: Google Alerts <googlealerts-noreply@google.com>

Subject: Google Alert - "John R. Frank"

==== Web - 2 new results for **["John R. Frank"]** ===

John R. Frank

SPOKANE, Wash. - **John R. Frank, 55**, died March 4, 2012, in **Coeur d' Alene, Idaho**. Survivors include: his wife, Miki; daughter, Patricia Frank; ...

<<http://www.hutchnews.com/obituaries/Frank--John-CP>>

In Memory of **John R Frank**

Biography. **John R. Frank, age 55**, passed away at **Sacred Heart Medical Center in Spokane, WA**, on March 4, 2012. **John** was born in Hutchison, KS, ...

<<http://www.englishfuneralchapel.com/sitemaker/sites/Englis1/obit.cgi?user=583335F>>

Key Questions

- What makes an alert “good”?
- How to measure it?
- How to scale entity-centric search?

==== Web - 2 new results for ["**John R. Frank**"] ===

John R. Frank

SPOKANE, Wash. - **John R. Frank**, 55, died March 4, 2012, in **Coeur d' Alene, Idaho**. Survivors include: his wife, Miki; daughter, Patricia Frank; ...

<<http://www.hutchnews.com/obituaries/Frank--John-CP>>

In Memory of **John R Frank**

Biography. **John R. Frank**, age 55, passed away at **Sacred Heart Medical Center in Spokane, WA**, on March 4, 2012. **John** was born in Hutchison, KS, ...
<<http://www.englishfuneralchapel.com/sitemaker/sites/Englis1/obit.cgi?user=583335F>>

Entities in Wikipedia or another Knowledge Base

especially if potentially libelous or harmful. (April 2008)

Takashi Murakami (村上 隆 Murakami Takashi[?], born in Tokyo) is an internationally prolific contemporary Japanese artist. He works in fine arts media—such as painting and sculpture—as well as what is conventionally considered commercial media—fashion, merchandise, and animation—and is known for blurring the line between high and low art. He coined the term superflat, which has become a buzzword in art circles.



Automatically recommend new edits



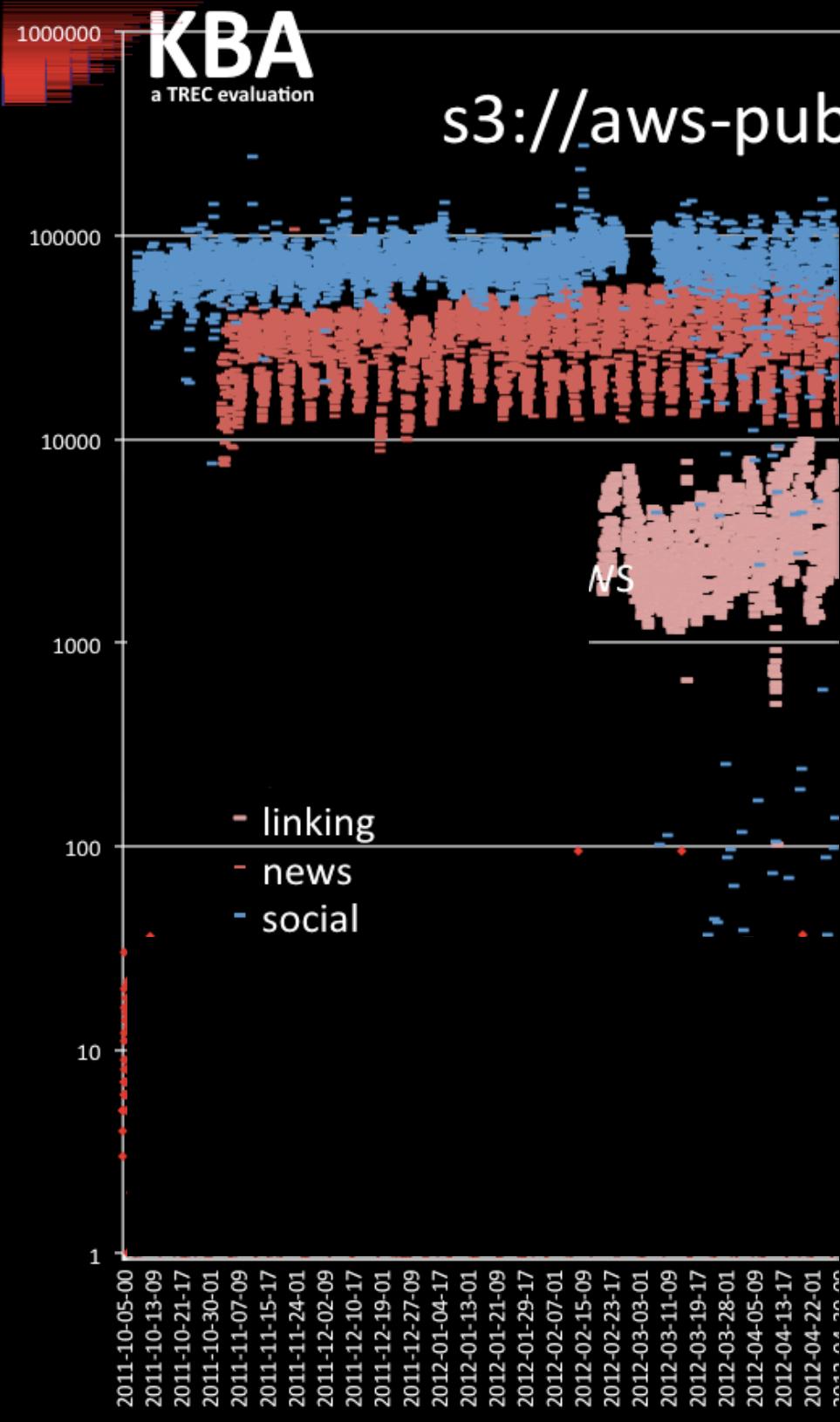
TREC KBA 2013

Stream Filtering to Recommend Citations and Fill Slots

- 1) Initialize with target entities
 - Start with profiles from Feb 2012
 - Training labels from Oct 2011–Feb 2012
- 2) Iterate over stream of text items
- 3) Identify “Vital” texts that change the profile
- 4) Identify slot values that change

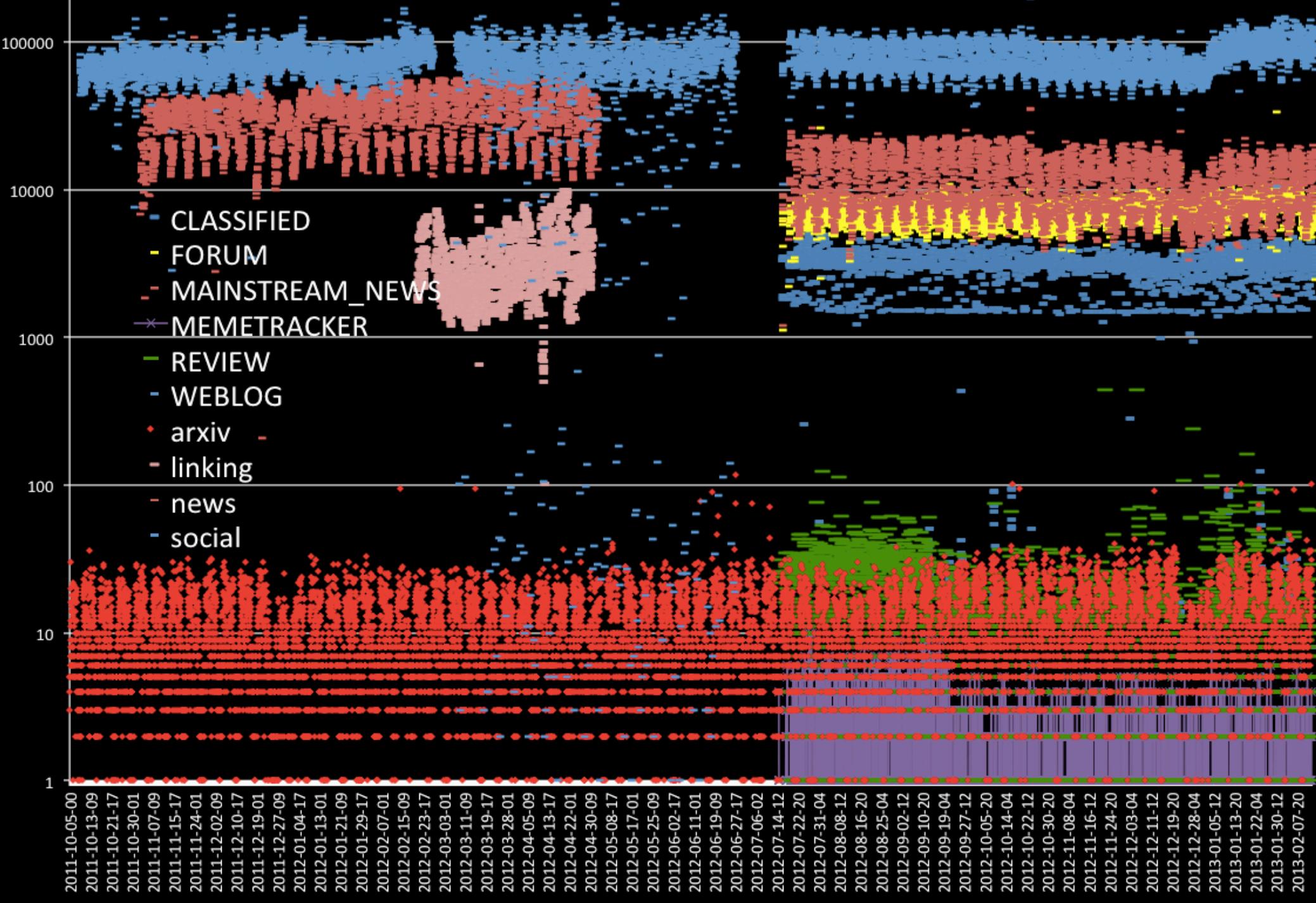
Content Stream

- 1bn texts, 50% English
- 11,948 hours of data (17months)
- 10^5 docs-per-hour
- News, blogs, forums, and link shortening





s3://aws-publicdatasets/trec/index.html



training evaluation

Corpus counts per hour

Pre-hoc assess all surface mentions

“vital” docs change entity profile; usually filling a slot

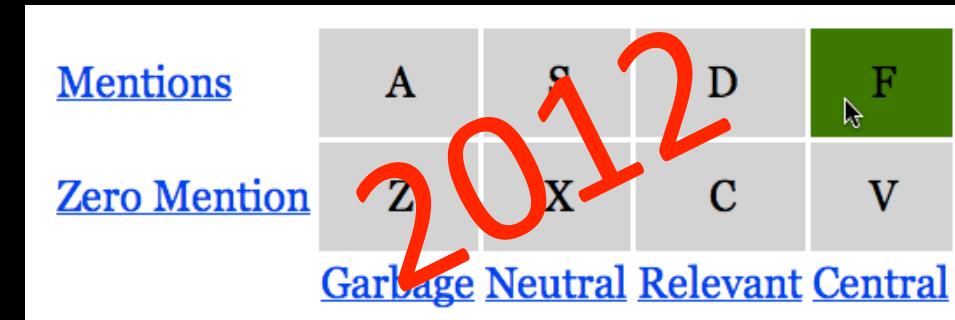
2011-10-05-00
2011-10-16-09
2011-10-22-01
2011-10-27-17
2011-11-02-09
2011-11-08-01
2011-11-13-17
2011-11-19-09
2011-12-05-01
2011-12-23-09
2011-12-29-01
2012-01-03-17
2012-01-09-09
2012-01-15-01
2012-01-20-17
2012-01-26-09
2012-02-01-01
2012-02-06-17
2012-02-12-09
2012-02-18-01
2012-02-23-17
2012-03-01-01
2012-03-11-17
2012-03-17-09
2012-03-23-01
2012-03-28-17
2012-04-03-09
2012-04-10-09
2012-04-20-09
2012-04-26-01
2012-05-01-17
2012-05-07-09
2012-05-13-01
2012-05-18-17
2012-05-24-09
2012-05-30-01
2012-06-04-17
2012-06-10-09
2012-06-16-01
2012-06-21-17
2012-06-27-09
2012-07-03-02
2012-07-08-20
2012-07-14-12
2012-07-20-04
2012-07-25-20
2012-08-01-12
2012-08-17-12
2012-08-23-04
2012-08-28-20
2012-09-03-12
2012-09-09-04
2012-09-14-20
2012-09-20-12
2012-09-26-04
2012-10-01-20
2012-10-07-12
2012-10-13-04
2012-10-18-20
2012-10-24-12
2012-10-30-04
2012-11-04-20
2012-11-10-12
2012-11-16-04
2012-11-21-20
2012-11-27-12
2012-12-03-04
2012-12-08-20
2012-12-14-12
2012-12-20-04
2012-12-25-20
2013-01-03-12
2013-01-06-04
2013-01-11-20
2013-01-17-12
2013-01-28-20
2013-02-03-12
2013-02-09-04

Refine “central” → “vital”

Masaru Emoto

From Wikipedia, the free encyclopedia

Masaru Emoto (江本 勝 *Emoto*
Masaru?, born July 22, 1943) is a
Japanese author and entrepreneur,



Published: March 31, 2012

Impact of Thoughts on Water

By Denis Gorce-Bourge

Water covers 70% of our Blue planet and our body is made of about 70% water.

Masaru Emoto is a Japanese Photographer and scientist. He is known over the world for his remarkable work on water and its deep connection with individual and collective consciousness.

For decades, Masaru took pictures of frozen crystals of water and tested the direct influence of the environment on the quality of those crystals.

Refine “central” → “vital”

Masaru Emoto

From Wikipedia, the free encyclopedia

Masaru Emoto (江本 勝 *Emoto*
Masaru?, born July 22, 1943) is a
Japanese author and entrepreneur,



Mentions	A	S	D	F
Zero Mention	Z	X	C	V
	Garbage	Neutral	Useful	Vital

2013

Published: March 31, 2012

Impact of Thoughts on Water

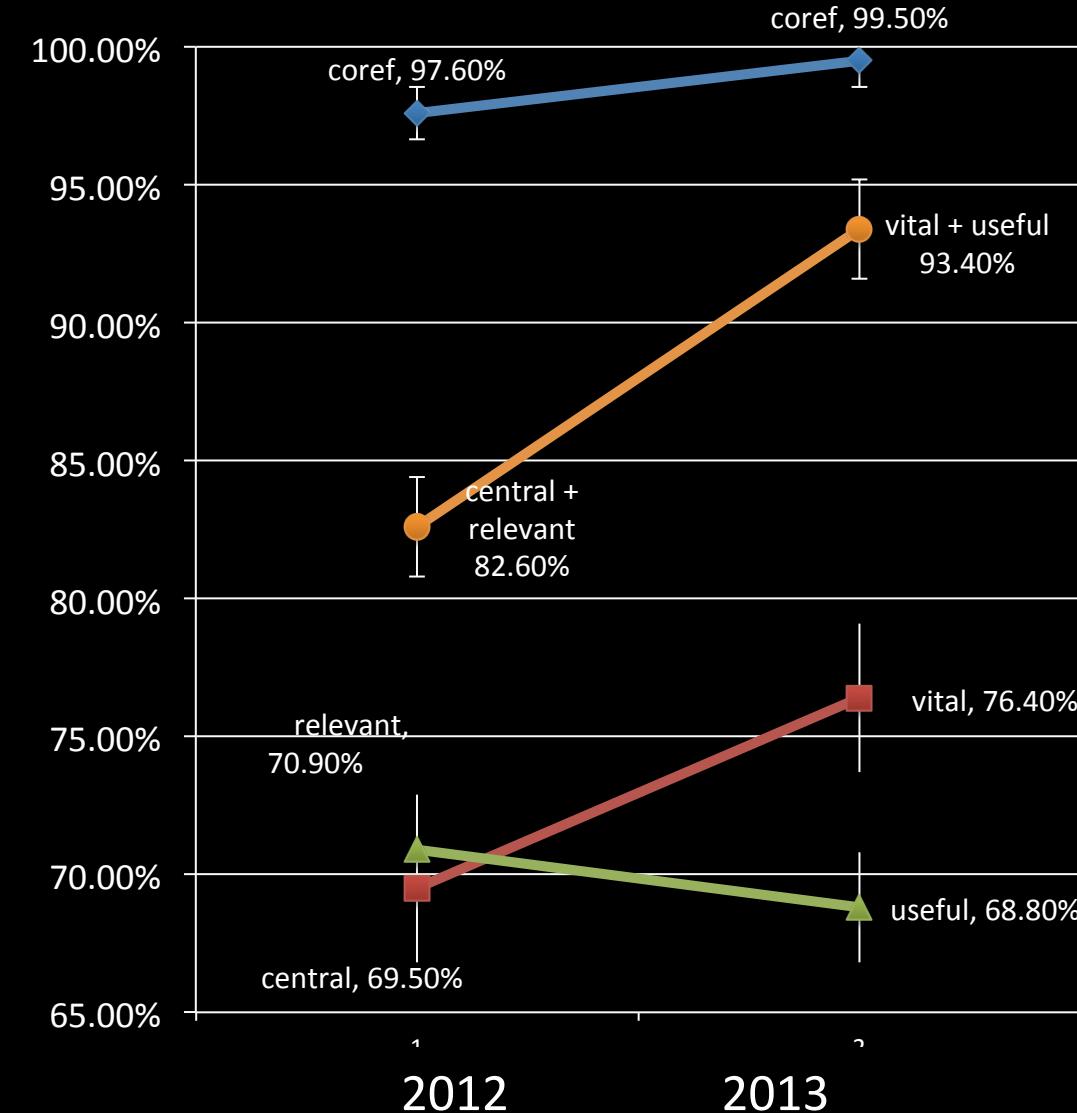
By Denis Gorce-Bourge

Water covers 70% of our Blue planet and our body is made of about 70% water.

Masaru Emoto is a Japanese Photographer and scientist. He is known over the world for his remarkable work on water and its deep connection with individual and collective consciousness.

For decades, Masaru took pictures of frozen crystals of water and tested the direct influence of the environment on the quality of those crystals.

Timeliness

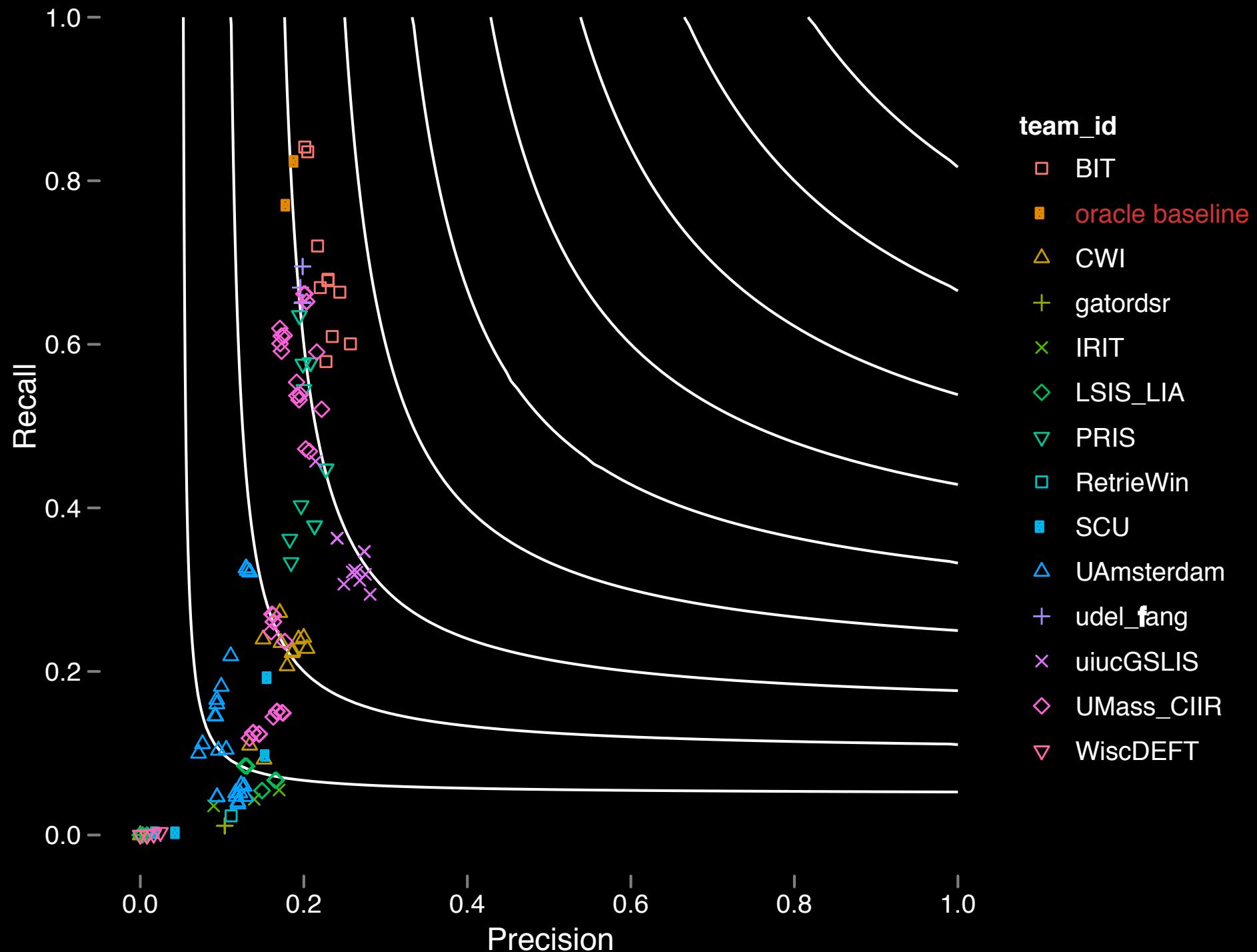


Mentions	A	S	D	F
Zero Mention	Z	X	C	V
	Garbage	Neutral	Useful	Vital

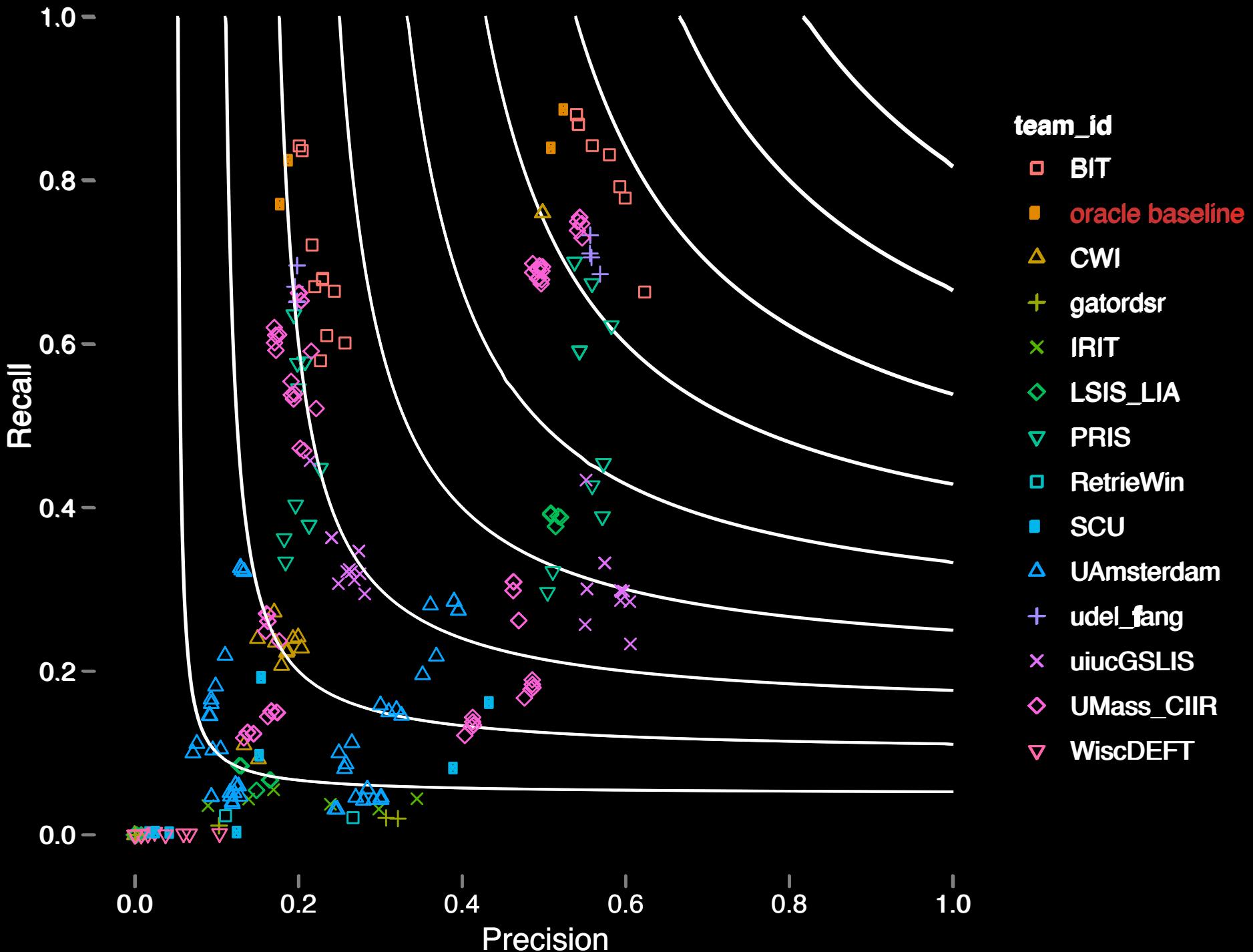
The majestic gray birds with red caps should be hundreds of miles south where it's warm at this time of year -- and not in Nebraska.
 "I've been there 50 years and I've never seen it," said noted ornithologist and author **Paul Johnsgard** of Lincoln.

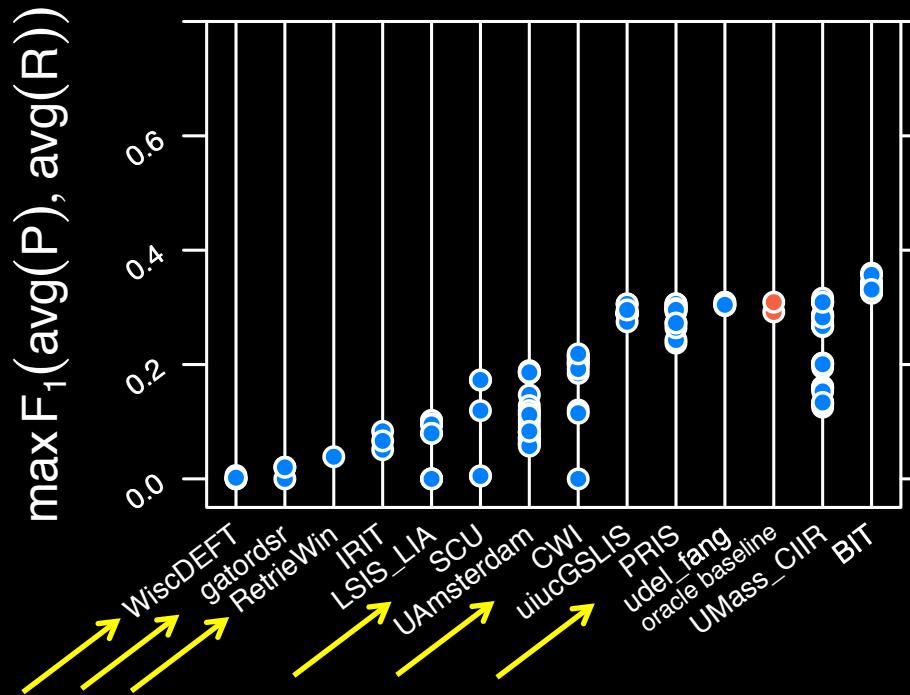
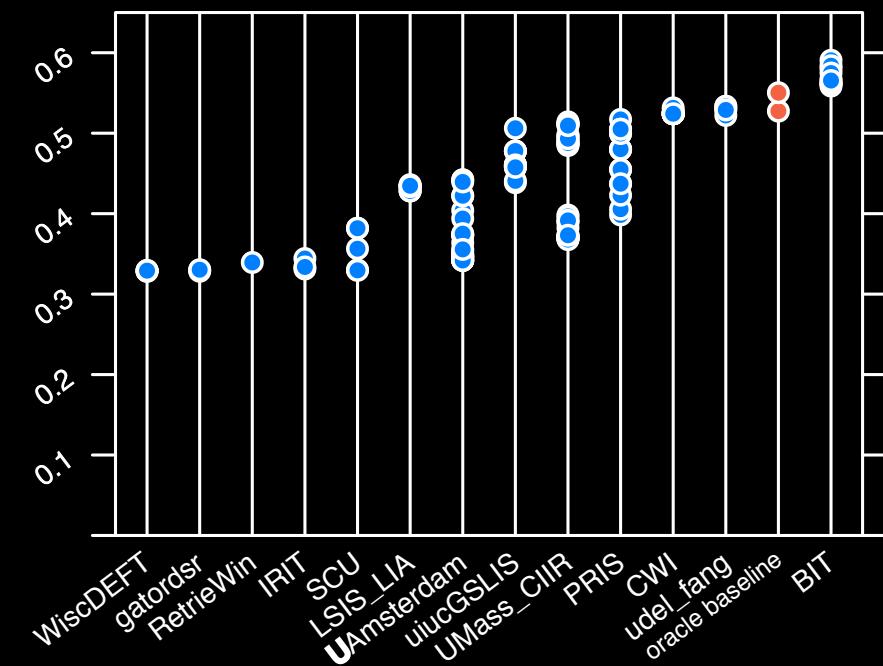
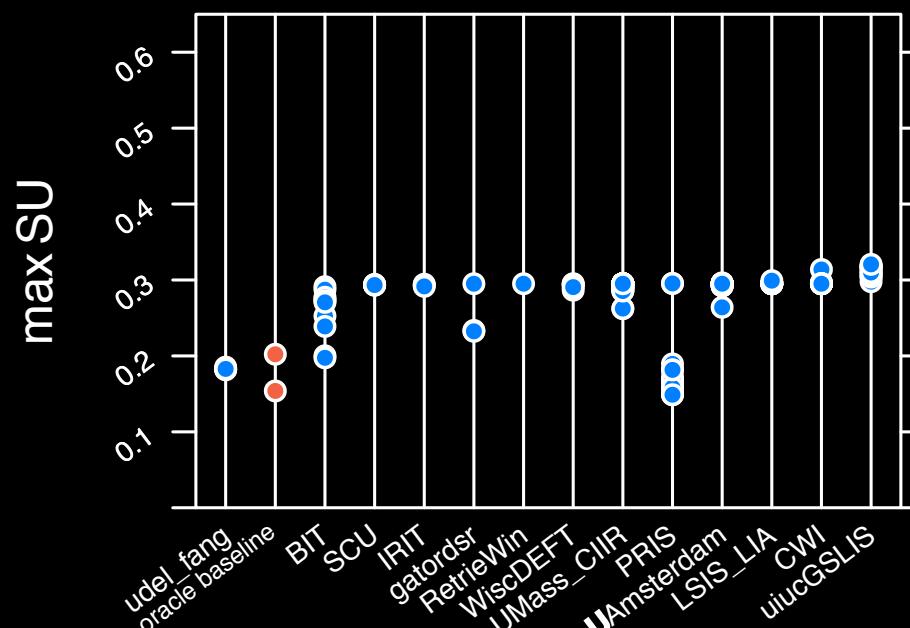
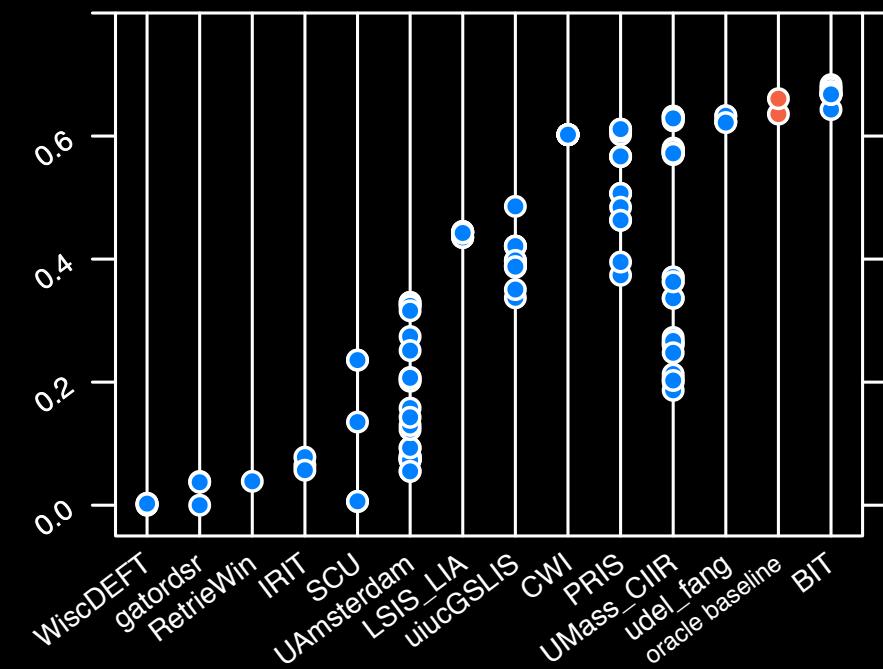


max F_1 (macroP, macroR)

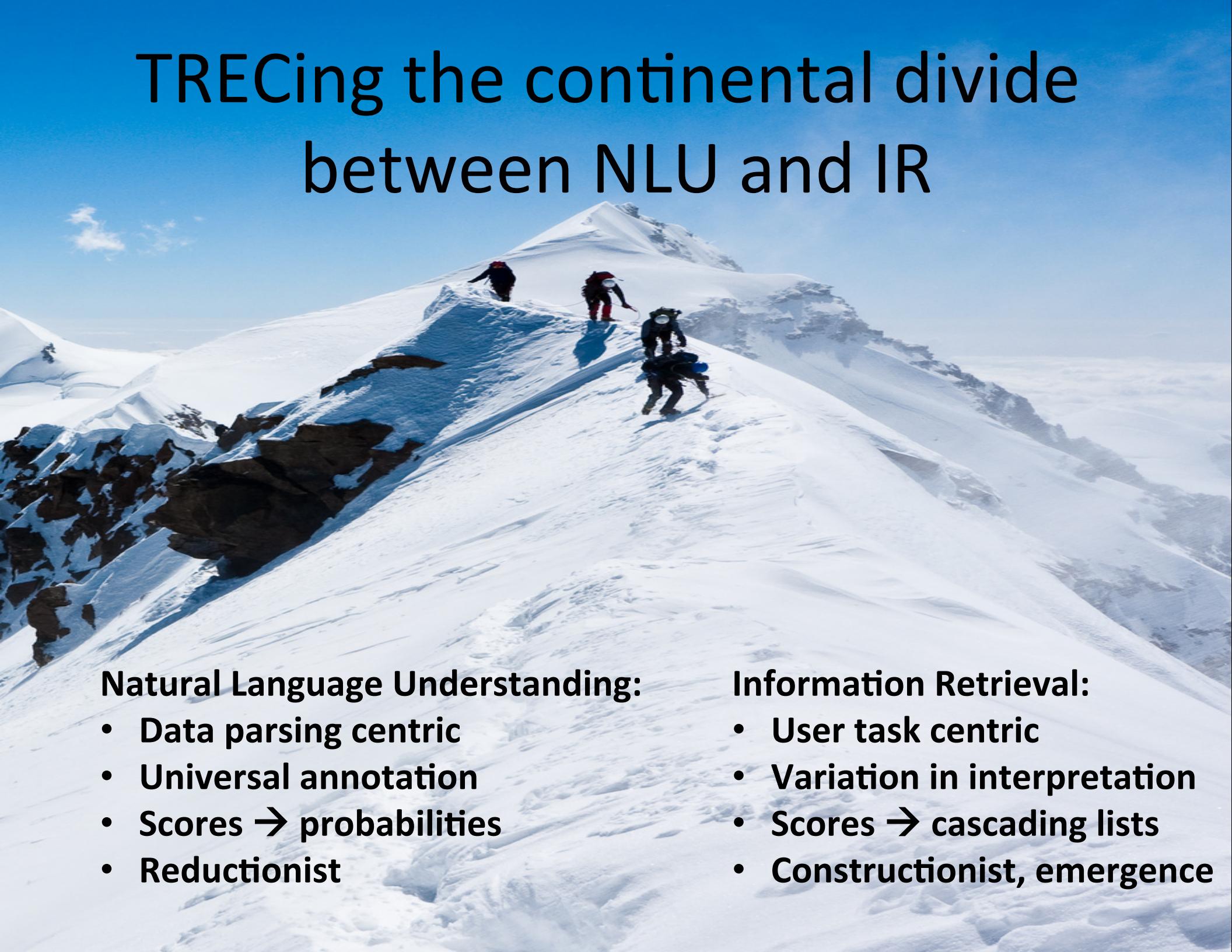


max F_1 (macroP, macroR)



vital

vital+useful


TRECing the continental divide between NLU and IR

A photograph of three climbers on a steep, snow-covered mountain ridge. They are wearing dark gear and are positioned at different points along the ridge, some higher up and some lower down. The background shows more of the mountain range under a clear blue sky.

Natural Language Understanding:

- Data parsing centric
- Universal annotation
- Scores → probabilities
- Reductionist

Information Retrieval:

- User task centric
- Variation in interpretation
- Scores → cascading lists
- Constructionist, emergence

Streaming Slot Filling

Person (PER)

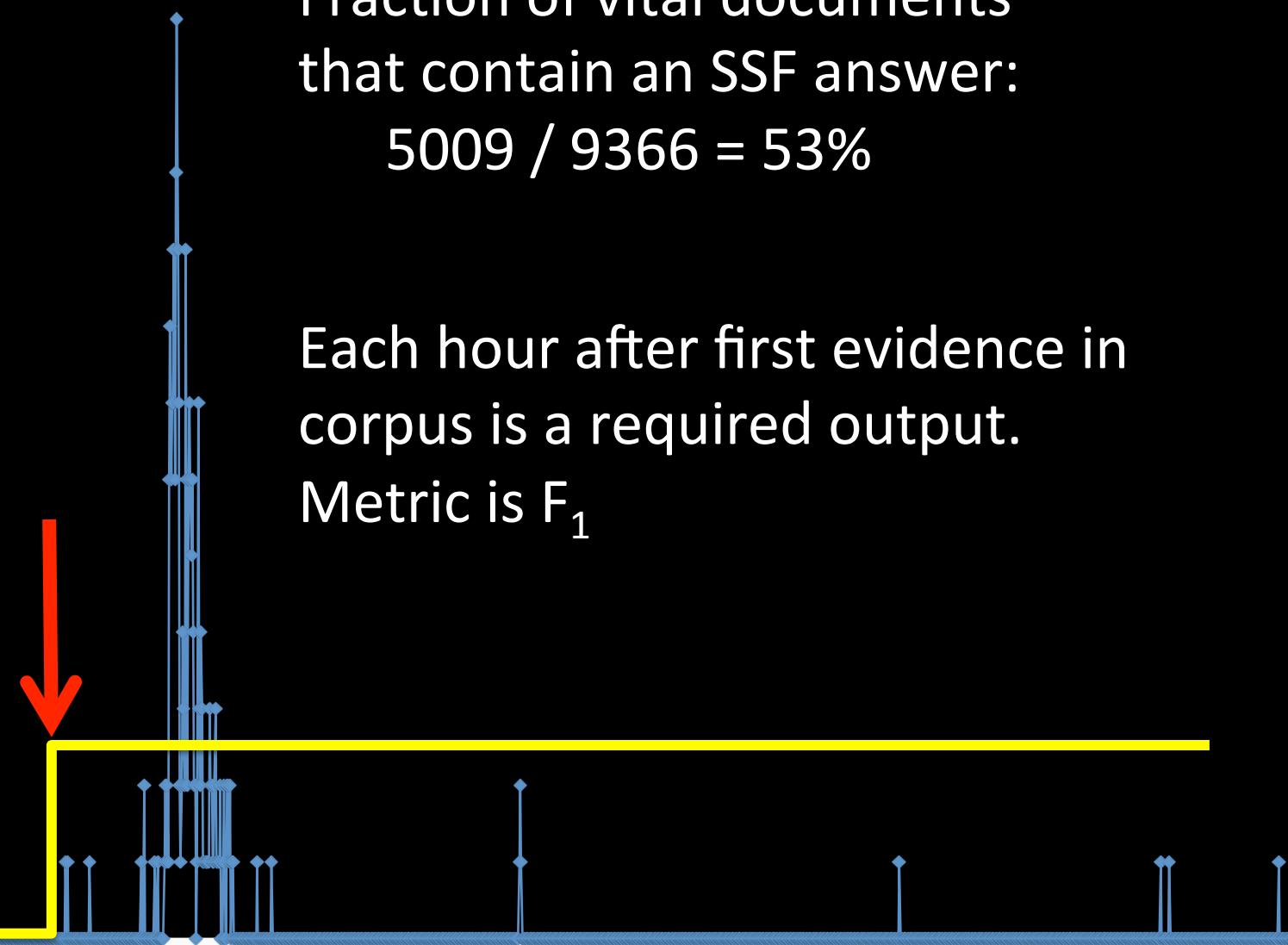
- Affiliate
- AssociateOf
- Contact_Meet_PlaceTime
- AwardsWon
- DateOfDeath
- CauseOfDeath
- Titles
- FounderOf
- EmployeeOf
- Significant Other
- Children

Facility (FAC)

- Affiliate
- Contact_Meet_Entity

Organization (ORG)

- Affiliate
- TopMembers
- FoundedBy

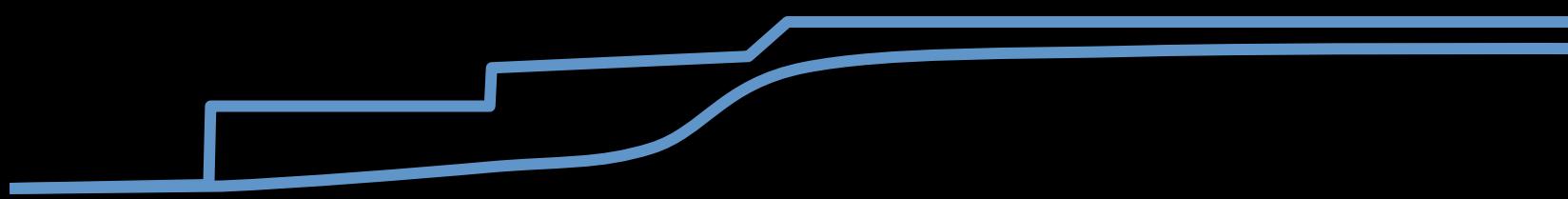
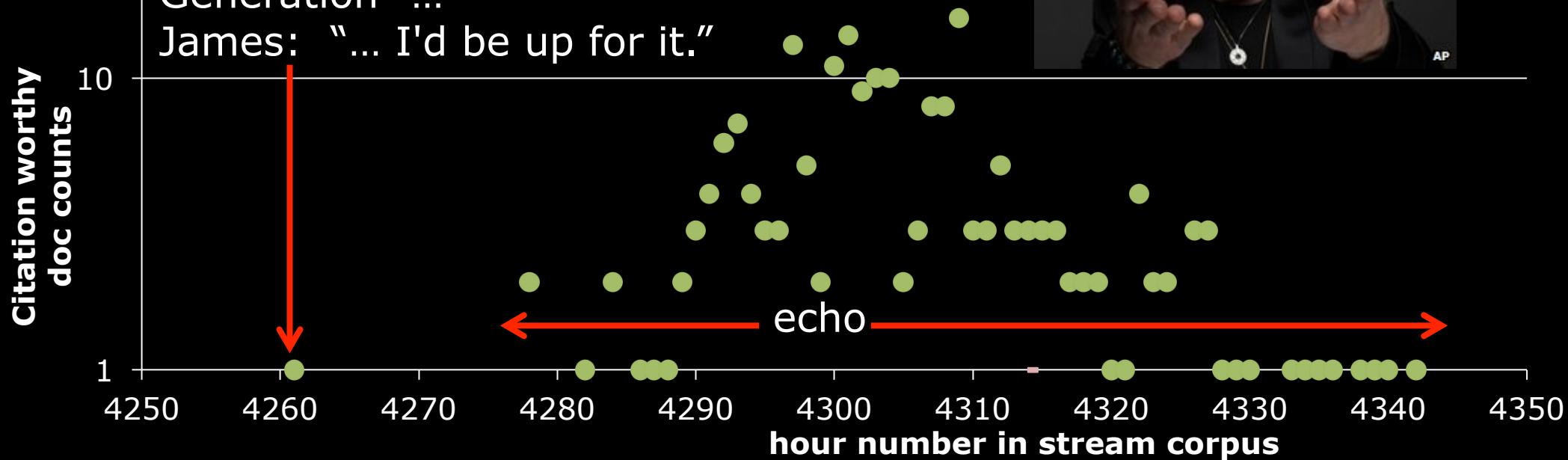


SSF example:

“FounderOf” slot on James McCartney

BBC: “What would you say to forming The Beatles - The Next Generation” ...

James: “... I'd be up for it.”

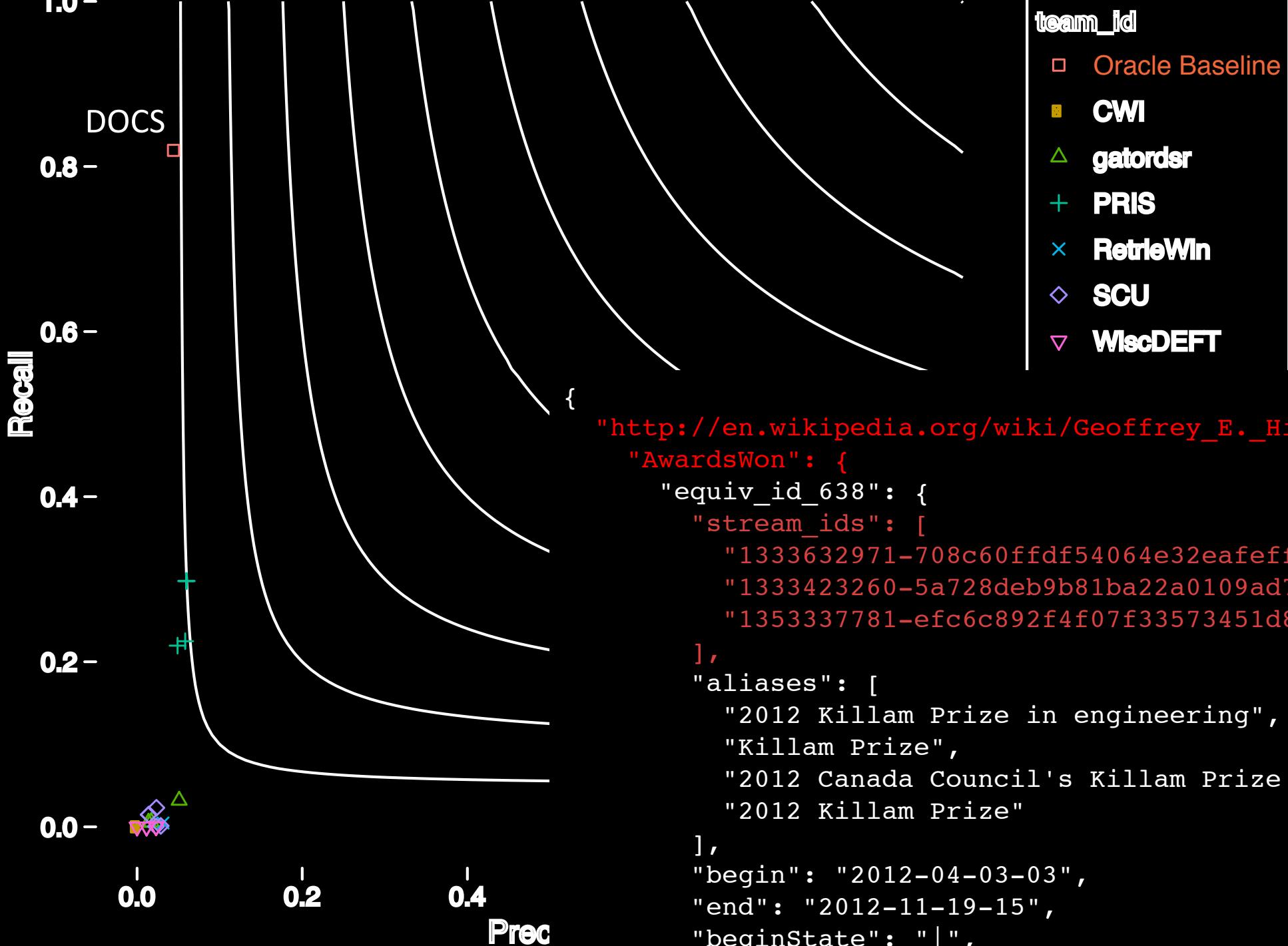


Confidence of sensitive versus insensitive systems

SSF Truth Record

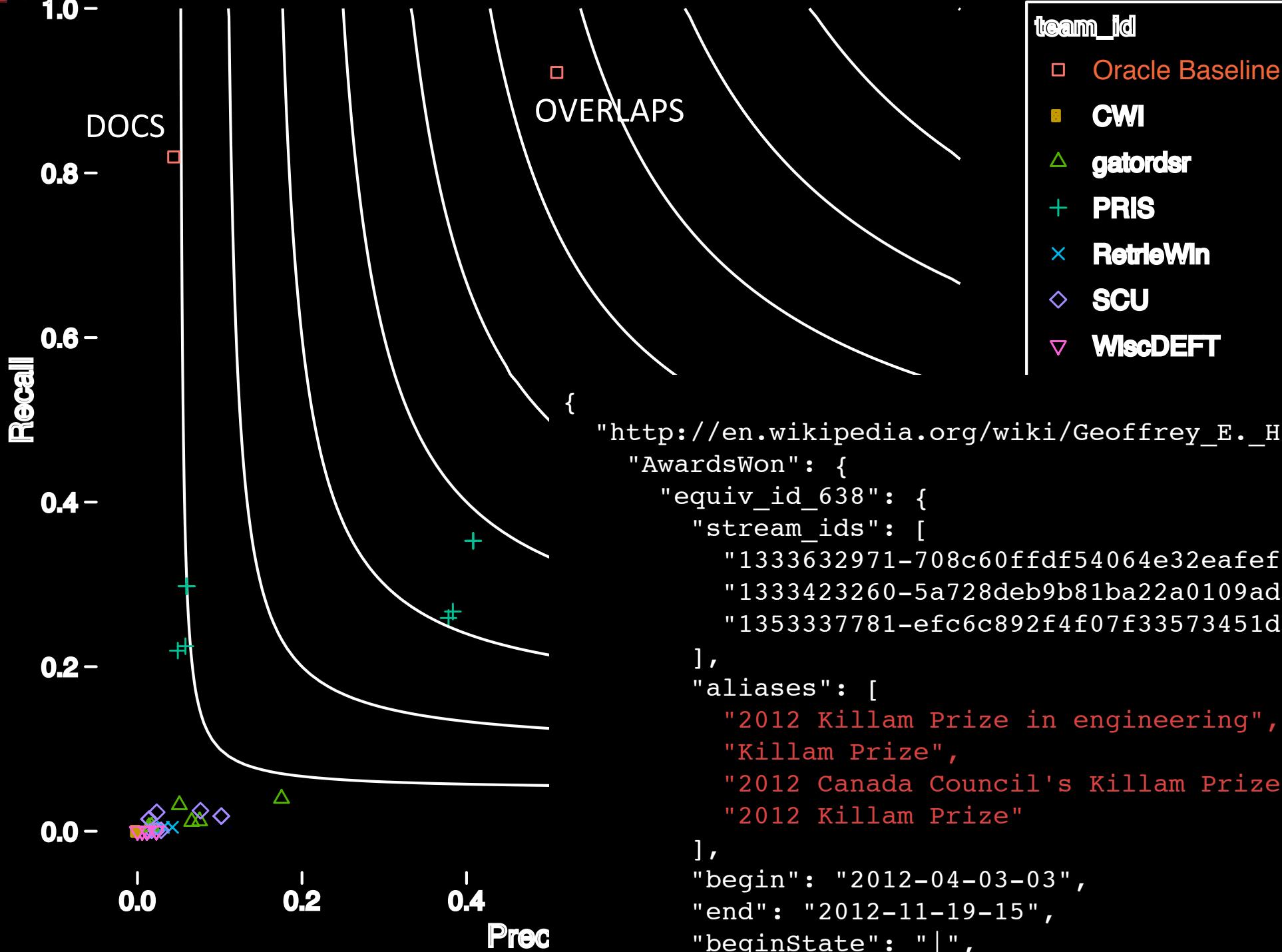
```
{  
  "http://en.wikipedia.org/wiki/Geoffrey_E._Hinton": {  
    "AwardsWon": {  
      "equiv_id_638": {  
        "stream_ids": [  
          "1333632971-708c60ffd54064e32eafeffbe145a51",  
          "1333423260-5a728deb9b81ba22a0109ad767ba4457",  
          "1353337781-efc6c892f4f07f33573451d818e8629c"  
        ],  
        "aliases": [  
          "2012 Killam Prize in engineering",  
          "Killam Prize",  
          "2012 Canada Council's Killam Prize laureates",  
          "2012 Killam Prize"  
        ],  
        "begin": "2012-04-03-03",  
        "end": "2012-11-19-15",  
        "beginState": "|",  
        "endState": ">"  
      },  
    }  
  }  
}
```

DOCS-level

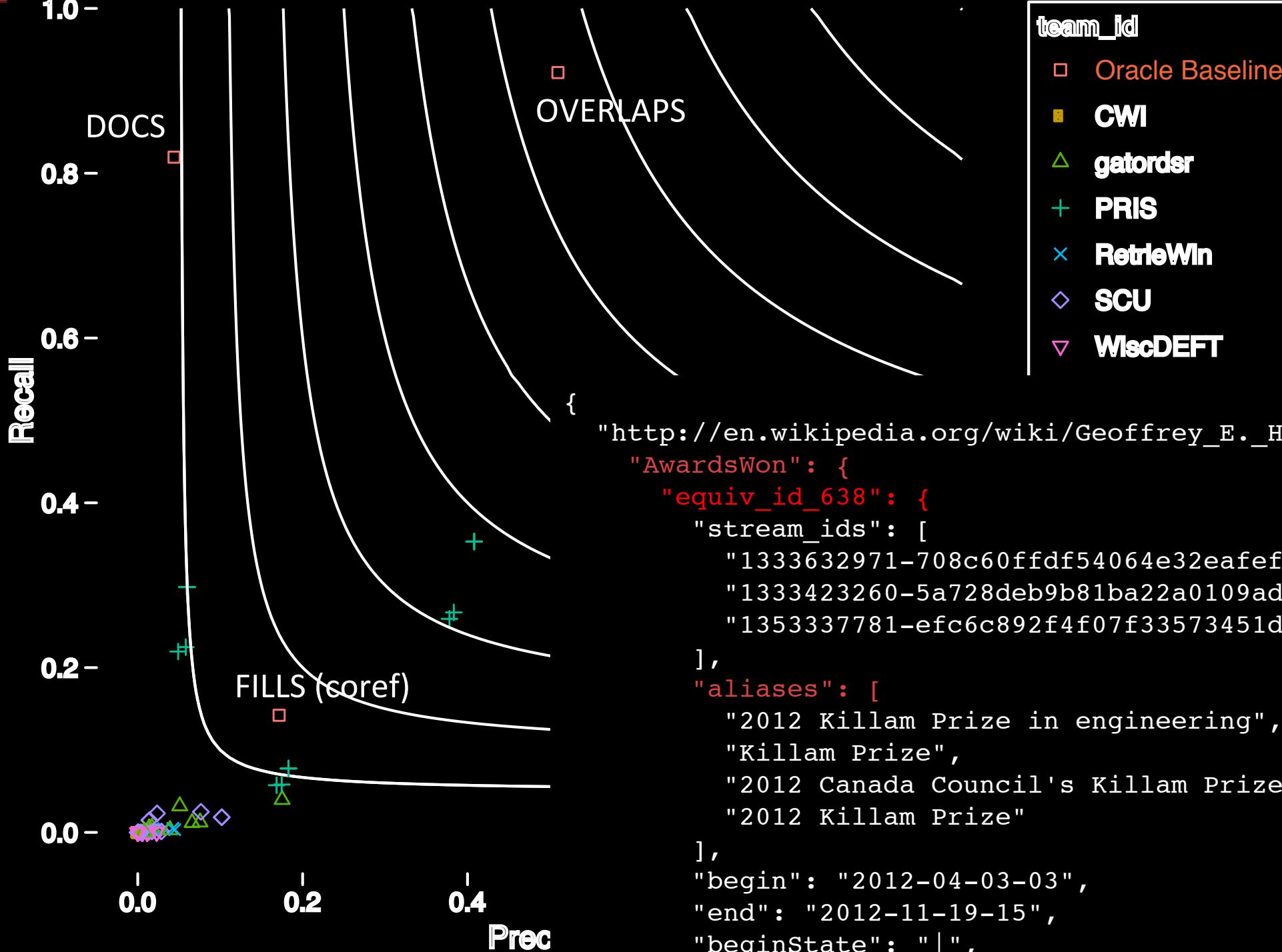


"http://en.wikipedia.org/wiki/Geoffrey_E._Hinton"
"AwardsWon": {
 "equiv_id_638": {
 "stream_ids": [
 "1333632971-708c60ffdf54064e32eafeffbe145",
 "1333423260-5a728deb9b81ba22a0109ad767ba4",
 "1353337781-efc6c892f4f07f33573451d818e86"
],
 "aliases": [
 "2012 Killam Prize in engineering",
 "Killam Prize",
 "2012 Canada Council's Killam Prize laureate",
 "2012 Killam Prize"
],
 "begin": "2012-04-03-03",
 "end": "2012-11-19-15",
 "beginState": "|",
 "endState": ">"
 }
}

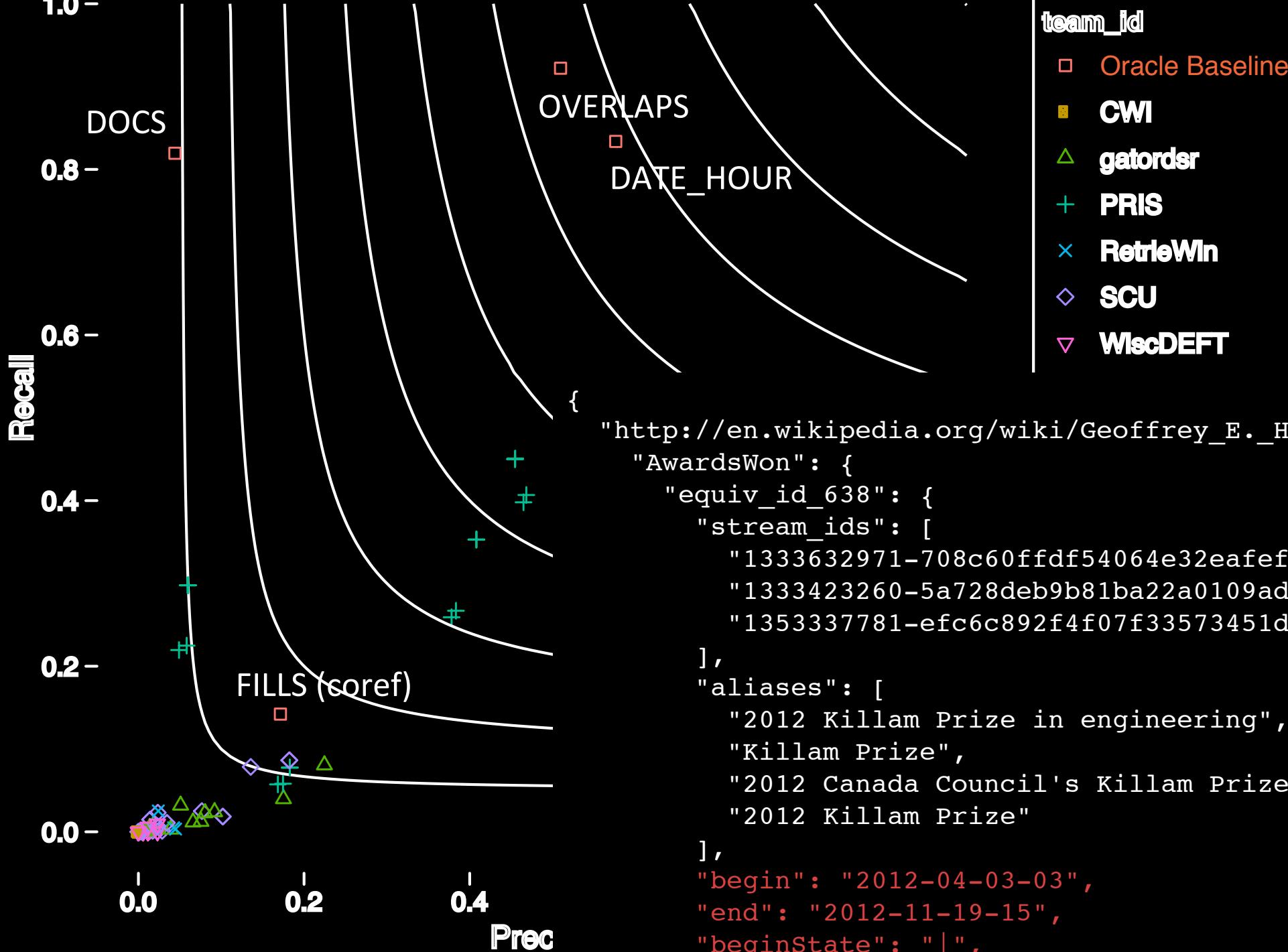
OVERLAPS-level



FILLS-level (coref)



DATE_HOUR-level



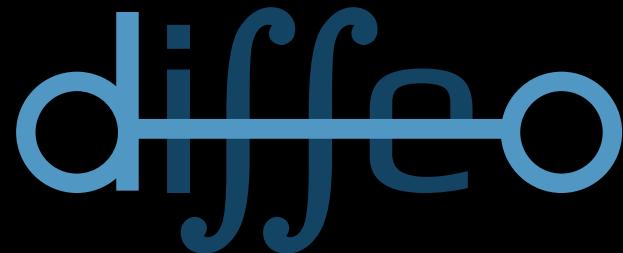
```

"http://en.wikipedia.org/wiki/Geoffrey_E._Hinton"
"AwardsWon": [
  "equiv_id_638": {
    "stream_ids": [
      "1333632971-708c60ffdf54064e32eafeffbe145",
      "1333423260-5a728deb9b81ba22a0109ad767ba4",
      "1353337781-efc6c892f4f07f33573451d818e86"
    ],
    "aliases": [
      "2012 Killam Prize in engineering",
      "Killam Prize",
      "2012 Canada Council's Killam Prize laureate",
      "2012 Killam Prize"
    ],
    "begin": "2012-04-03-03",
    "end": "2012-11-19-15",
    "beginState": "|",
    "endState": ">"
  }
]
  
```

KBA 2014

- Repeat CCR
 - more training data?
 - open source fast name filter
 - <http://github.com/trec-kba/streamcorpus-filter>
- Integrate SSF more closely with CCR
 - Do we need slots that are not in ACE or TAC?
- New metrics
 - Time-aware or slot-aware?

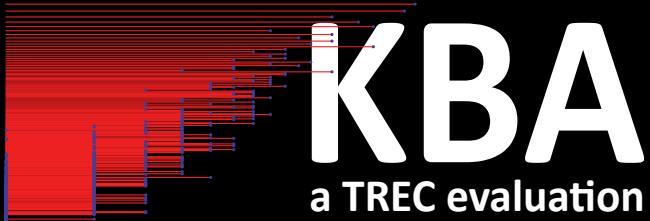
Sponsors: Thank You.



spinn3r

Raytheon
BBN Technologies





Thanks for your time.

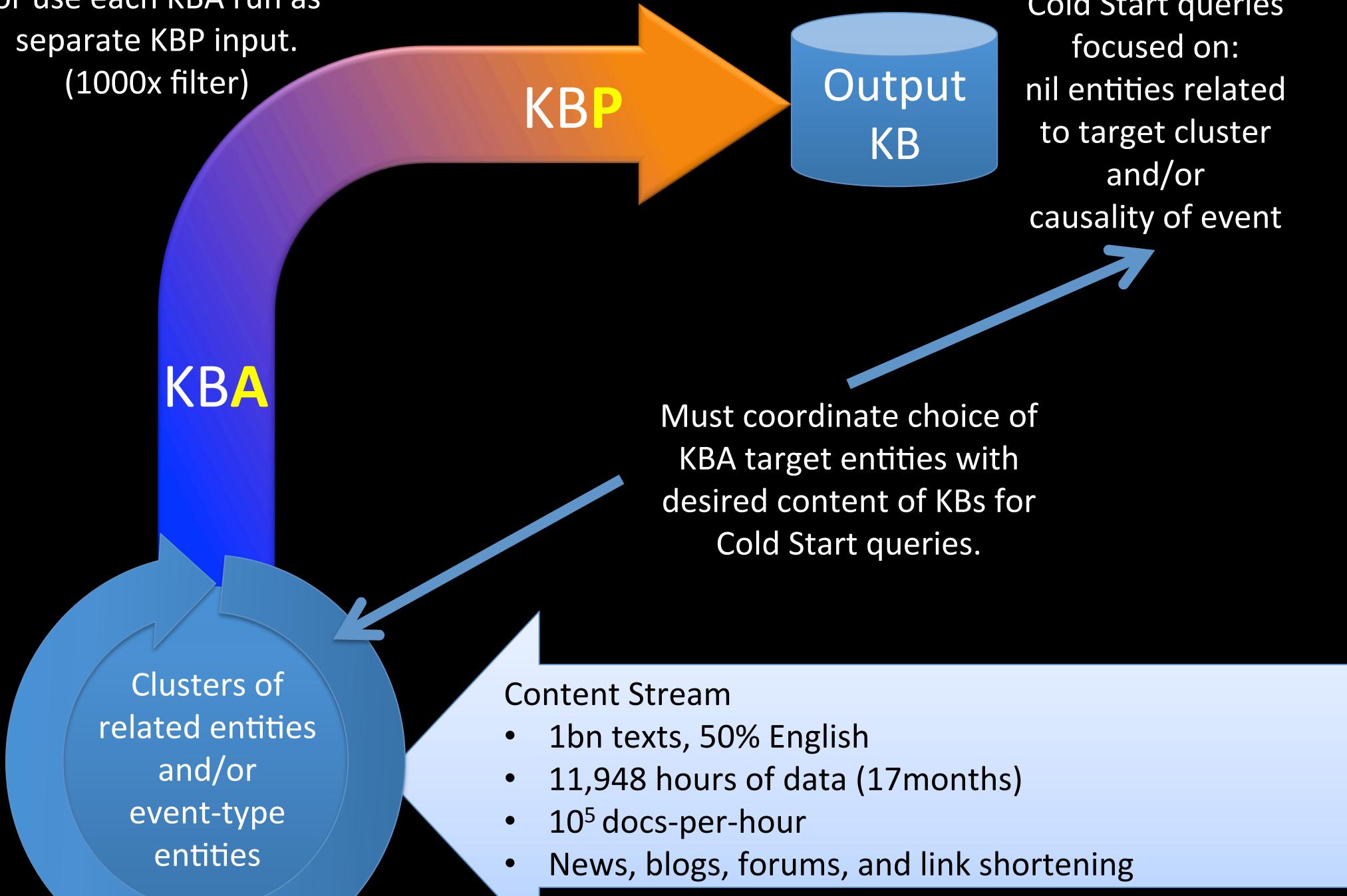
John R. Frank

jrf@mit.edu

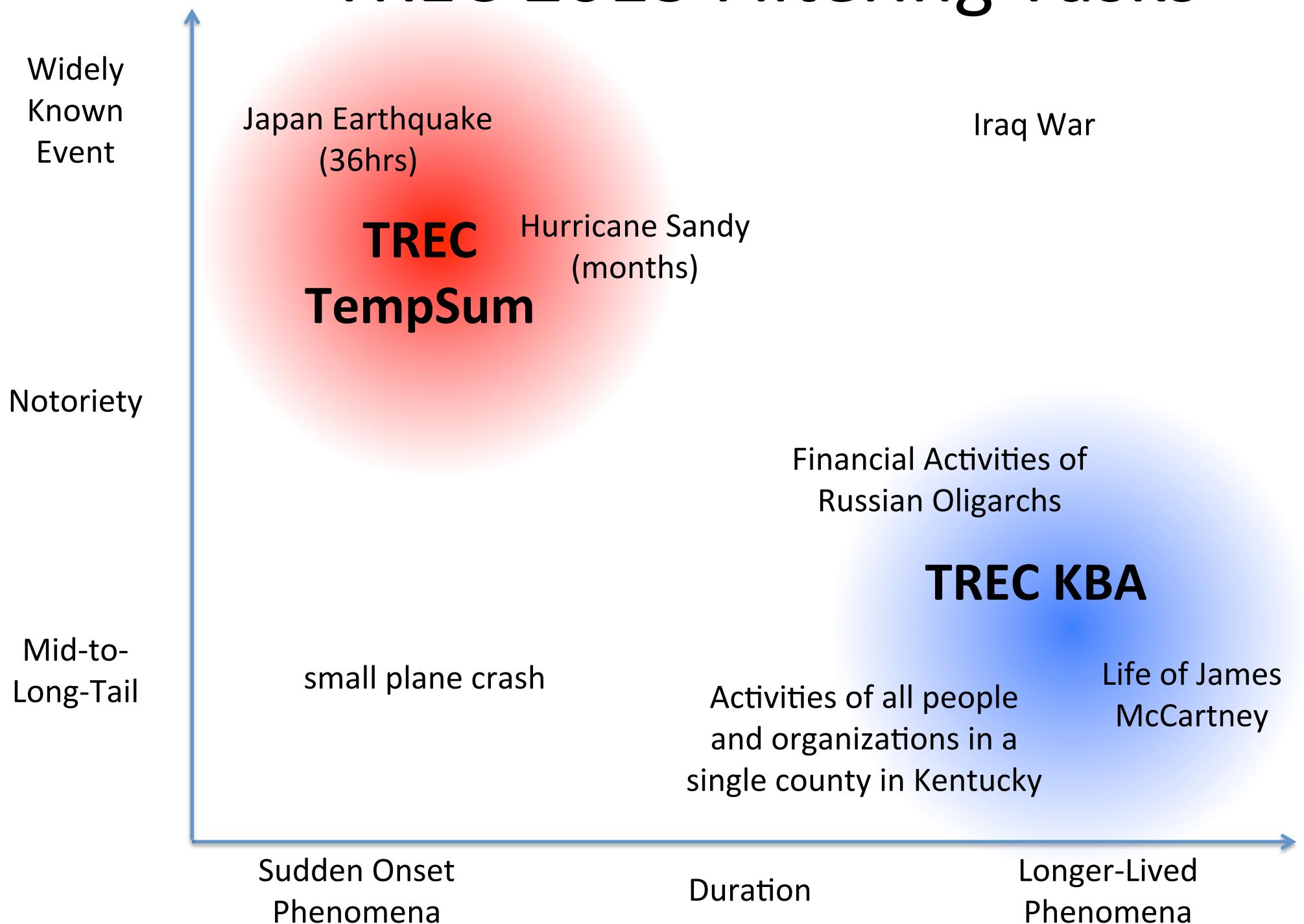
<http://trec-kba.org>

KBX

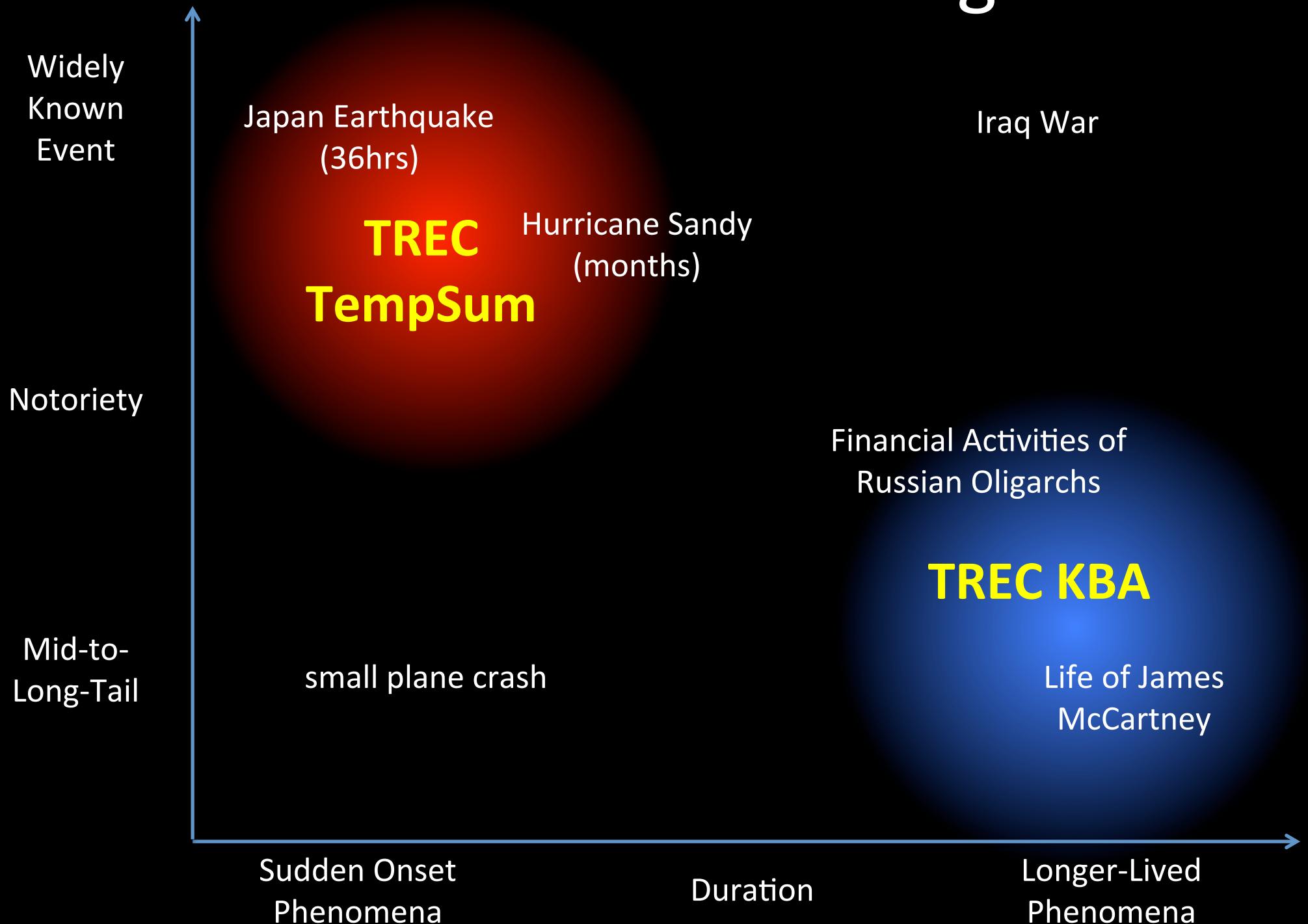
Pool top-K filtered docs,
or use each KBA run as
separate KBP input.
(1000x filter)



TREC 2013 Filtering Tasks

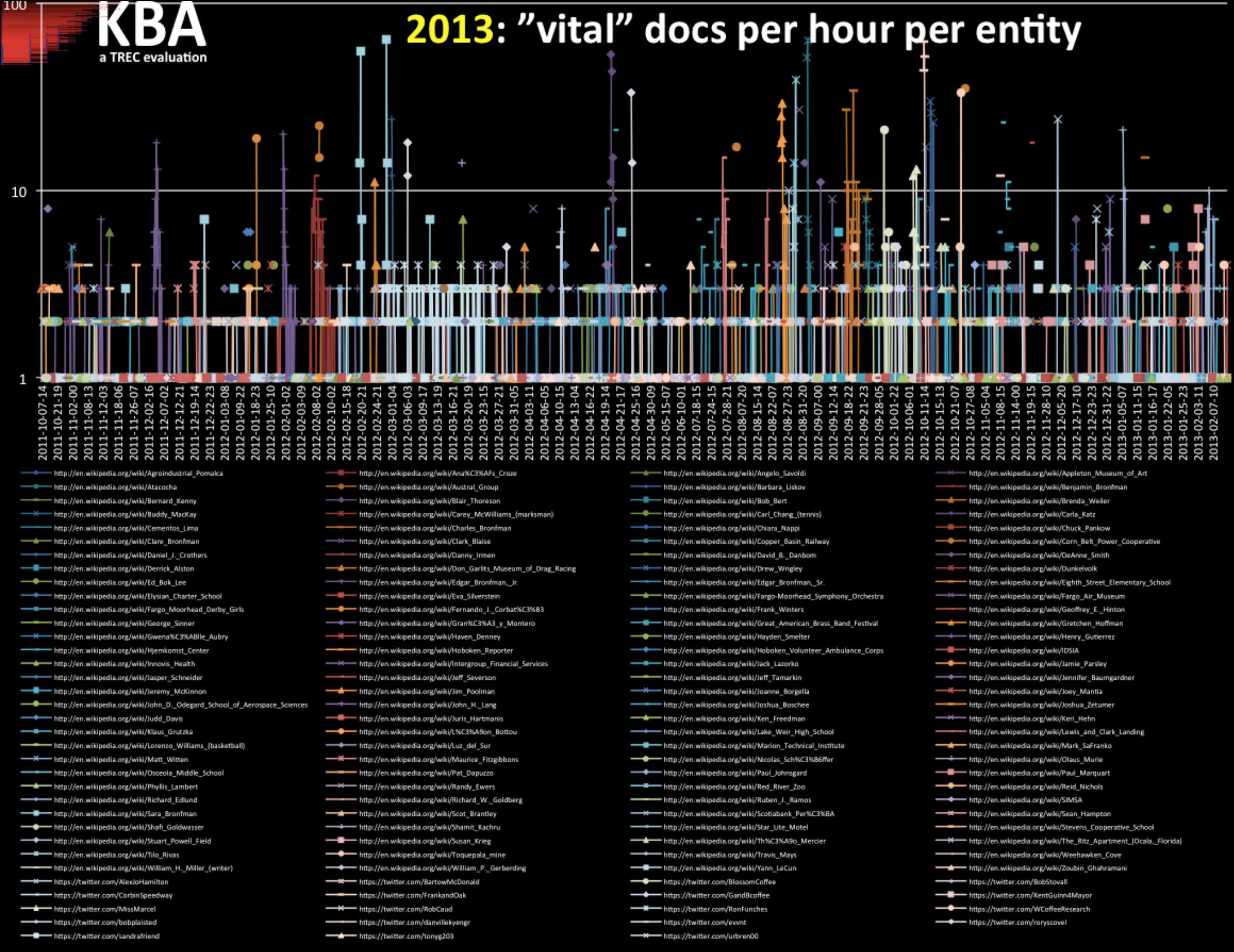


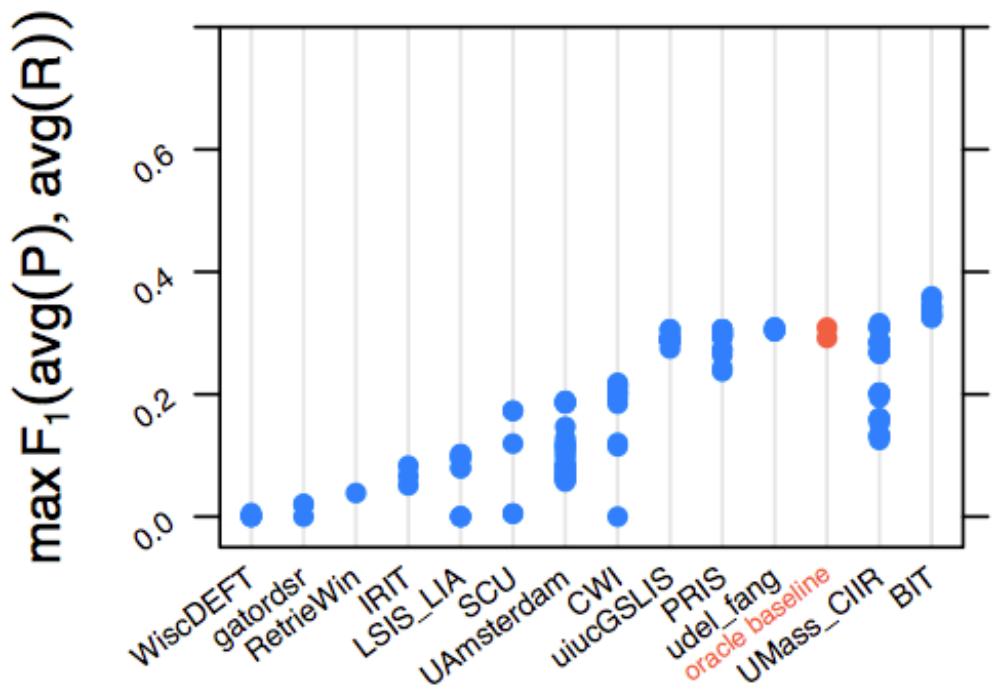
TREC 2013 Filtering Tasks



- Beijing Institute of Technology
- Centrum Wiskunde & Informatica
- Computable Insights
- Institt de Recherche en Informatique de Toulouse
- LSIS/LIA
- Pattern Recognition and Intelligent System laboratory
- Stanford University
- Santa Clara University
- University of Amsterdam
- InfoLab at University of Delaware
- University of Florida CISE Dept
- Grad. School of Lib. & Info Sci; Univ. of Illinois, Urbana-Champaign
- CIIR, School of Computer Science, Univ. of Massachusetts Amherst
- Wisconsin DEFT Research Group

2013: "vital" docs per hour per entity



vital**vital+useful**