# Learning to rank adaptively for scalable information extraction

**Pablo Barrio**, Columbia University
Gonçalo Simões, INESC-ID and IST, University of Lisbon
Helena Galhardas, INESC-ID and IST, University of Lisbon
Luis Gravano, Columbia University
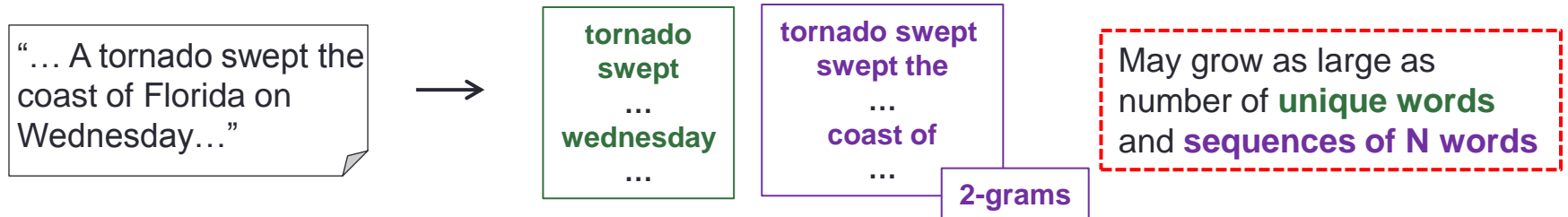
# Information Extraction (IE)

- Natural-language text **embeds** "structured" data

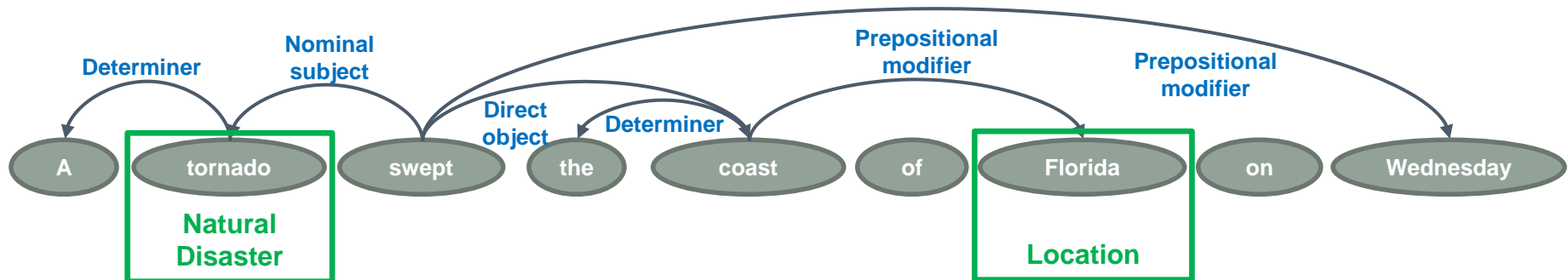- Information extraction systems **extract** this data

"… A **tornado** swept the coast of **Florida** on Wednesday…"

**Natural Disaster-Location**
information extraction system

<**tornado**, **Florida**>

Extracted tuple for
**Natural Disaster-Location**
relation

**Much richer querying and analysis possible**

# IE is Challenging and Time Consuming

- ## Operates over **large sets of features**
  Bag of words, N-grams, grammar productions, dependency paths

"… A tornado swept the coast of Florida on Wednesday…"

→

**tornado swept** ... **wednesday** ...

**tornado swept swept the** ... **coast of** ...

**2-grams**

May grow as large as number of **unique words** and **sequences of N words**

- ## Requires **complex text analysis**
  Dependency parsing, entity recognition, syntactic parsing, shallow parsing, part-of-speech tagging, semantic role labeling



May take **several seconds per document**
(e.g., with subsequence kernel extractor for Natural Disaster-Location)
Problematic over large document collections

# Reducing Processing Time: Opportunities

Documents are "useful" if they produce output for a given IE task

- **Small**, **topic-specific** fraction of collection

  Only **2% of documents** in a New York Times archive, mostly **environment-related**, are useful for Natural Disaster-Location with a state-of-the-art IE system

  > Should focus extraction over these documents and ignore rest

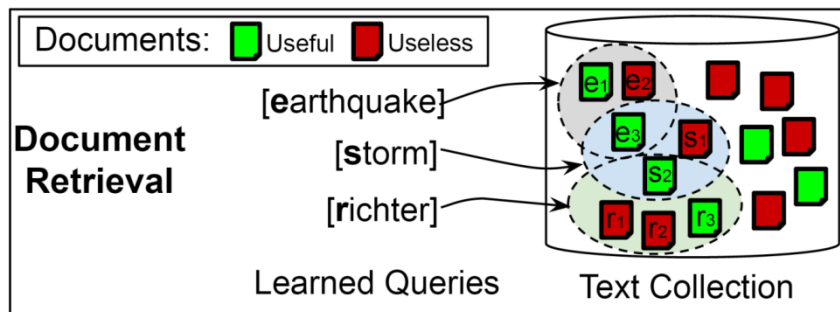- Useful documents share **distinctive words and phrases**

  "Earthquake," "storm," "Richter," "volcano eruption" for Natural Disaster-Location

  > Can learn to differentiate between useful documents for an IE task and rest
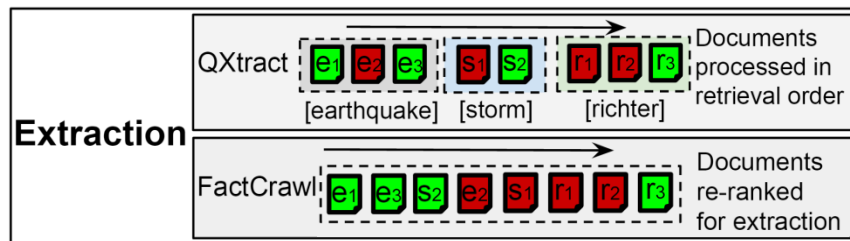
- Information extraction system **"labels"** documents as useful or not **for free**

  > IE process generates ever-expanding training set for learning to identify useful documents

# Existing Approaches: QXtract and FactCrawl



QXtract and FactCrawl learn from small document sample and exhibit far-from-perfect recall

FactCrawl ranks documents using learned queries and does not adapt to new processed documents
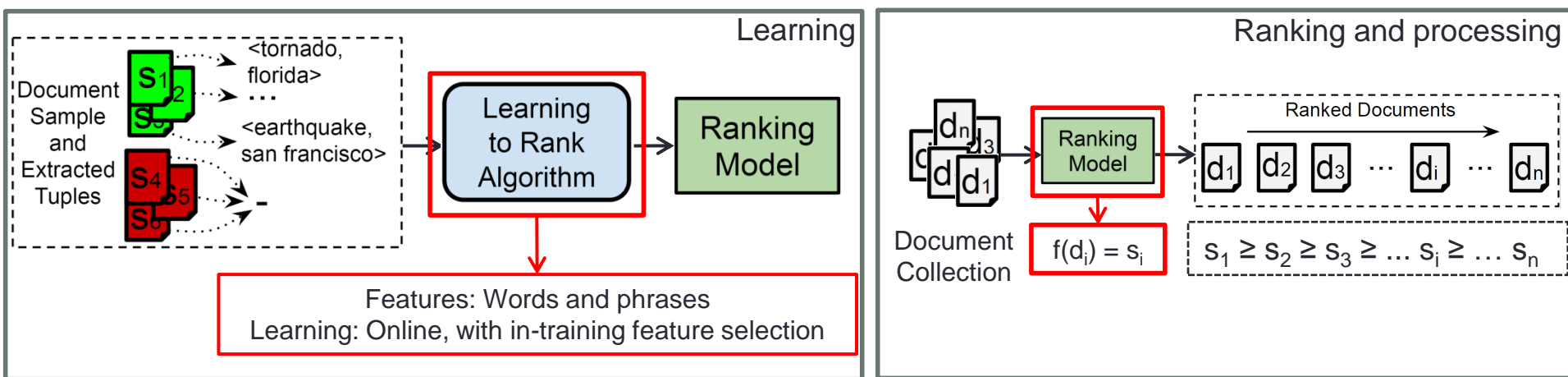
[Eugene Agichtein and Luis Gravano, "Querying text databases for efficient information extraction." *ICDE '03*]
[Christoph Boden et al., "FactCrawl: A fact retrieval framework for full-text indices." *WebDB '11*]
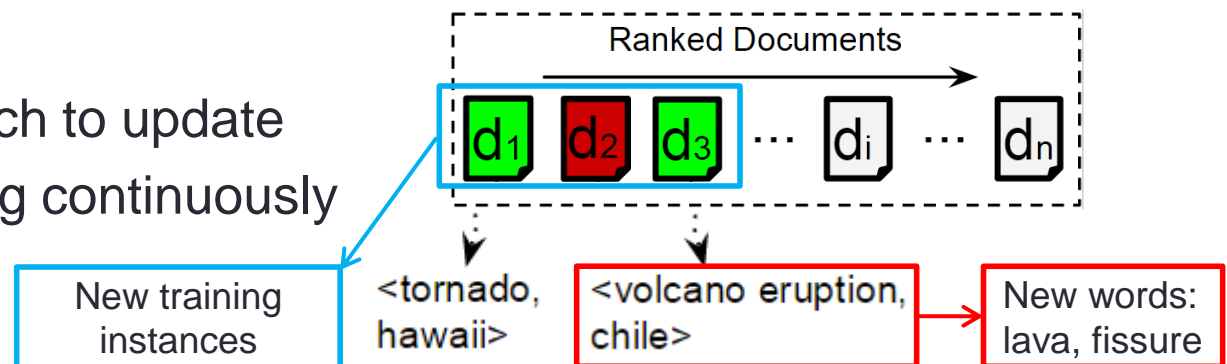
# Our Approach: Key Aspects

- Document ranking needs to be **robust and efficient**

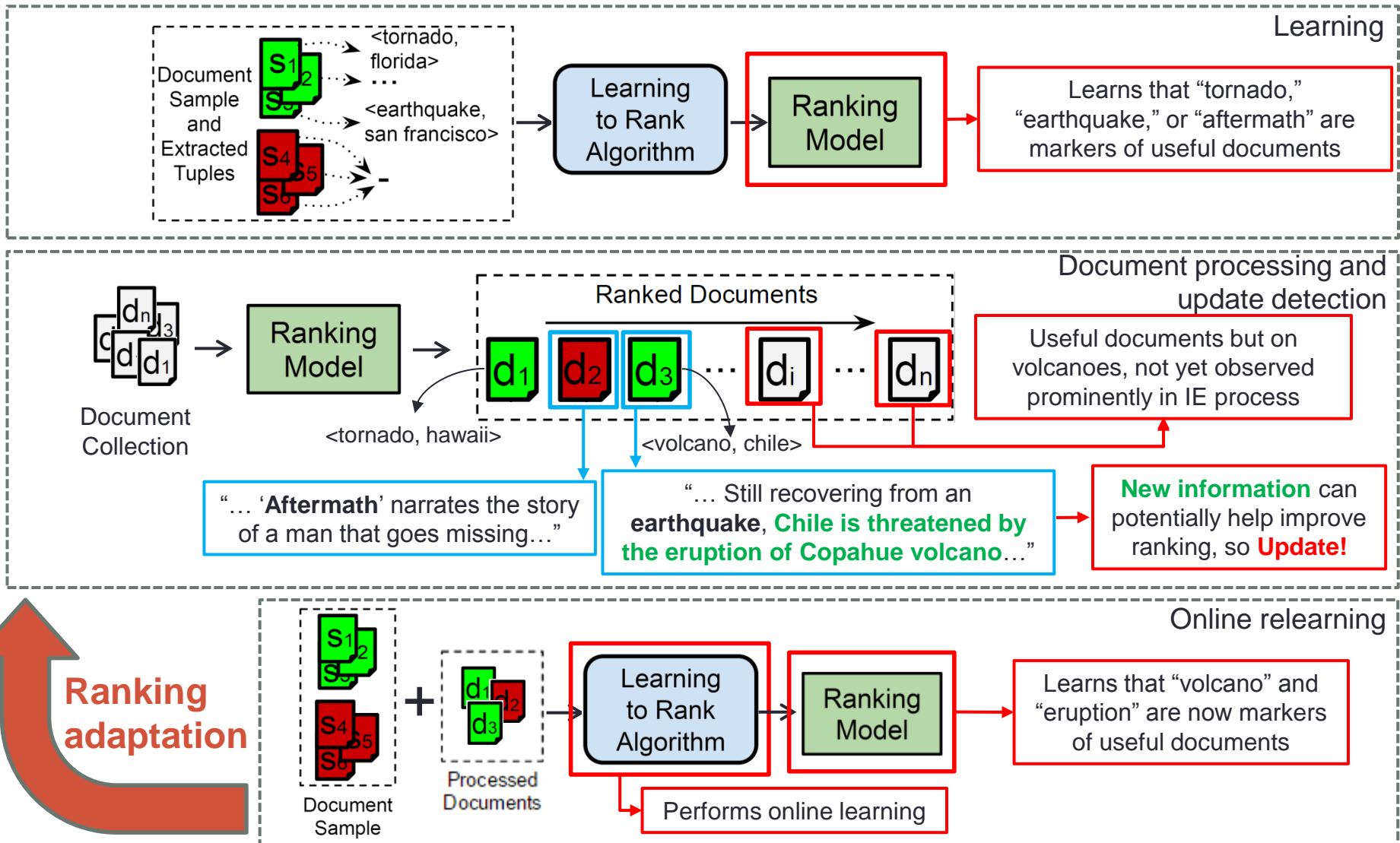  Learning to rank approach for document ranking



- Results of extraction process form **ever-expanding training set**

  Adaptive approach to update document ranking continuously
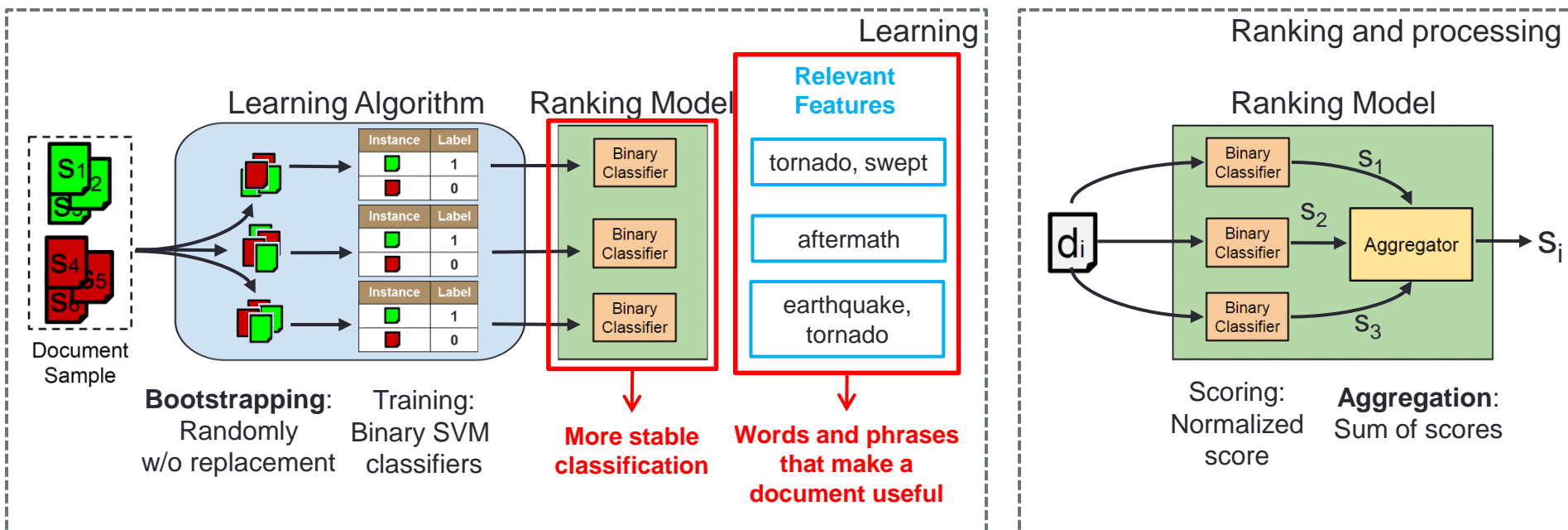
# Ranking Documents Adaptively for IE



**Learning**

Document Sample and Extracted Tuples → <tornado, florida> … <earthquake, san francisco> - → Learning to Rank Algorithm → Ranking Model → Learns that "tornado," "earthquake," or "aftermath" are markers of useful documents

**Document processing and update detection**

Document Collection → Ranking Model → Ranked Documents: $d_1$ $d_2$ $d_3$ … $d_i$ … $d_n$

<tornado, hawaii>
<volcano, chile>

"… 'Aftermath' narrates the story of a man that goes missing…"

"… Still recovering from an earthquake, Chile is threatened by the eruption of Copahue volcano…"

Useful documents but on volcanoes, not yet observed prominently in IE process

New information can potentially help improve ranking, so Update!

**Online relearning**

Ranking adaptation

Document Sample + Processed Documents → Learning to Rank Algorithm → Ranking Model → Learns that "volcano" and "eruption" are now markers of useful documents

Performs online learning

# Ranking Documents Adaptively for IE: Our Alternatives

- Efficient learning-to-rank techniques for information extraction: **BAgg-IE**, **RSVM-IE**

- Update detection techniques for document ranking adaptation: **Top-$K$**, **Mod-C**

# Efficient Learning to Rank for IE: BAgg-IE

- Based on **bootstrapping aggregation**



**Learning**

Document Sample

Learning Algorithm

| Instance | Label |
|----------|-------|
| | 1 |
| | 0 |

Ranking Model

Binary Classifier

**Relevant Features**

tornado, swept

aftermath

earthquake, tornado

**Bootstrapping**: Randomly w/o replacement

Training: Binary SVM classifiers

**More stable classification**

**Words and phrases that make a document useful**

**Ranking and processing**

Ranking Model

$d_i$ — Binary Classifier — $s_1$ — Aggregator — $s_i$

$s_2$

$s_3$

Scoring: Normalized score

**Aggregation**: Sum of scores

All models are trained using online learning and in-training feature selection
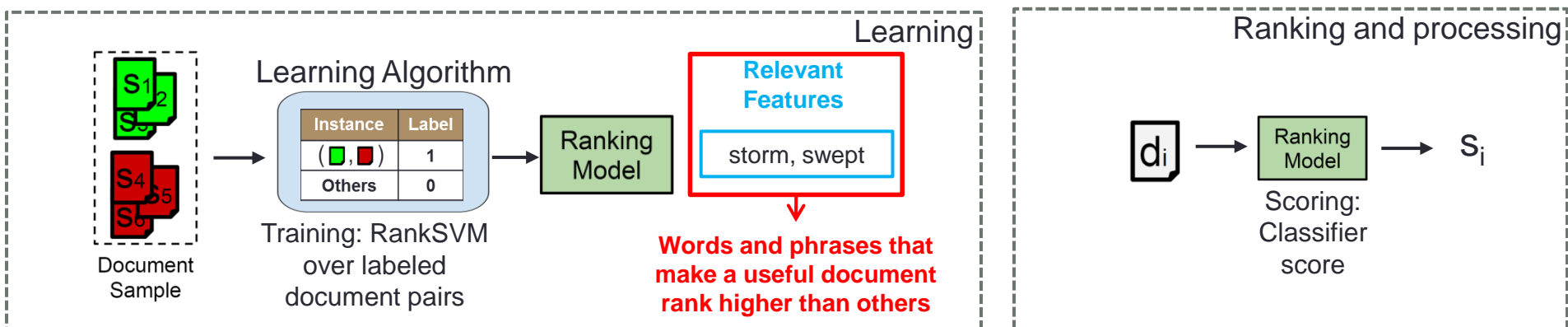
# Efficient Learning to Rank for IE: RSVM-IE

- Based on **RankSVM**

  Learns SVM classifier on pairwise
  difference of documents

| Learning model | Training instance |
|:---:|:---:|
| **RankSVM** | $d_i$ - $d_n$ |
| SVM | $d_i$ |

Training Label is **1** iif $d_i$ is "better" than $d_n$



Learning

Learning Algorithm

| Instance | Label |
|:---:|:---:|
| ( ■ , ■ ) | 1 |
| Others | 0 |

Training: RankSVM over labeled document pairs

Document Sample

Ranking Model

**Relevant Features**

storm, swept

**Words and phrases that make a useful document rank higher than others**

Ranking and processing

$d_i$ → Ranking Model → $S_i$

Scoring: Classifier score

Model is trained using online learning and in-training feature selection

# Ranking Documents Adaptively for IE: Our Alternatives

- Efficient learning-to-rank techniques for information extraction: **BAgg-IE**, **RSVM-IE**

- Update detection techniques for document ranking adaptation: **Top-*K***, **Mod-C**
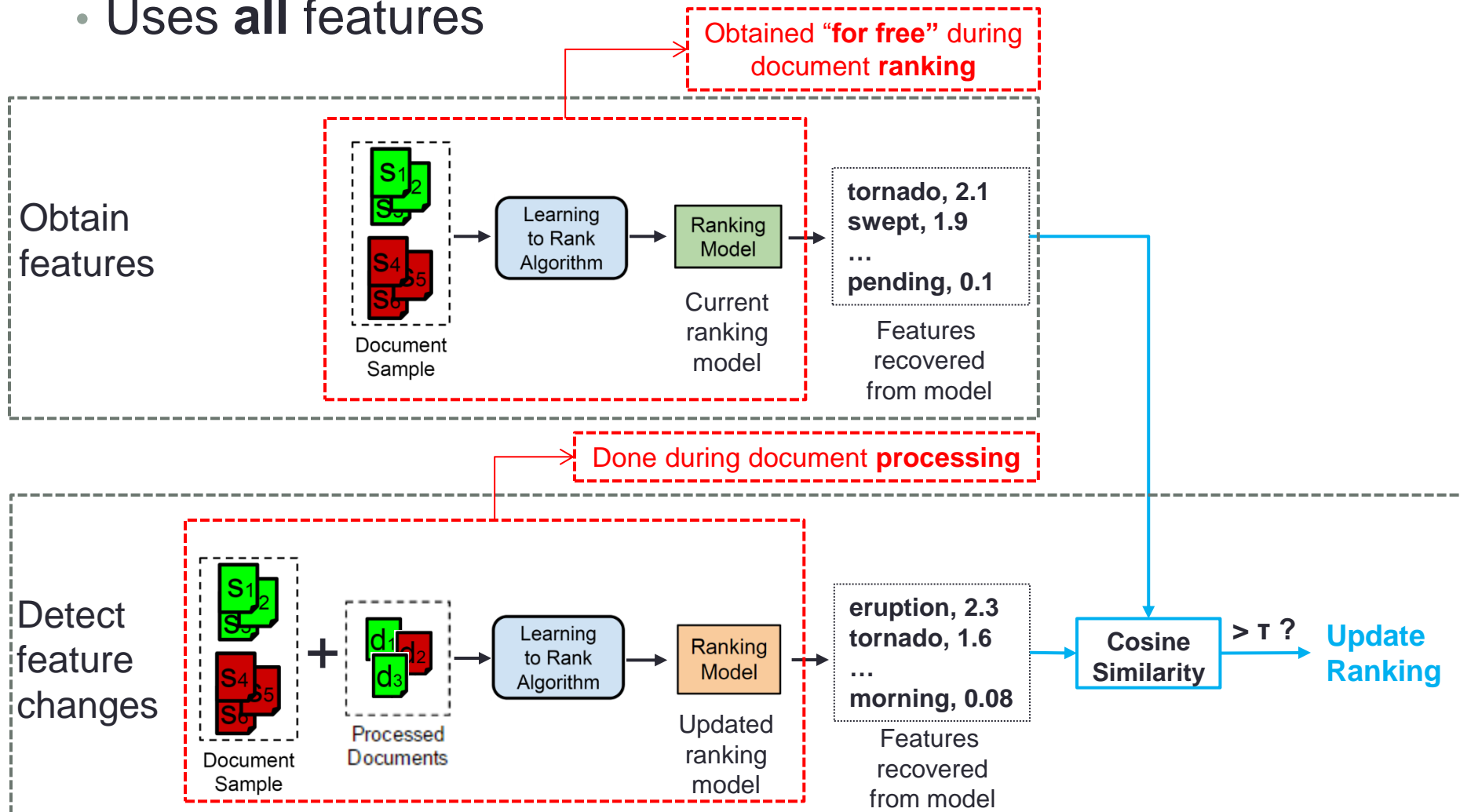
# Update Detection for Document Ranking Adaptation: Top-$K$

- Uses only **most important** (top-$K$) features

# Update Detection for Document Ranking Adaptation: Mod-C

- Uses **all** features



Obtained "**for free**" during document **ranking**

**Obtain features**

S$_1$ S$_2$ S$_3$
S$_4$ S$_5$ S$_6$

Document Sample

Learning to Rank Algorithm

Ranking Model

Current ranking model

**tornado, 2.1**
**swept, 1.9**
**…**
**pending, 0.1**

Features recovered from model

Done during document **processing**

**Detect feature changes**

S$_1$ S$_2$ S$_3$
S$_4$ S$_5$ S$_6$

Document Sample

+

d$_1$ d$_2$
d$_3$

Processed Documents

Learning to Rank Algorithm

Ranking Model

Updated ranking model

**eruption, 2.3**
**tornado, 1.6**
**…**
**morning, 0.08**

Features recovered from model

**Cosine Similarity**

**> τ ?**

**Update Ranking**

# Experimental Settings

- Dataset: **The New York Times** archive: 1.8 million articles from 1987-2007
- Information extraction systems

Simple extraction systems:
HMMs, text patterns

### Person-Organization

**Google** co-founders **Larry Page** and **Sergey Brin** recently sat down with billionaire venture capitalist Vinod Khosla for a lengthy interview.

| Person | Organization |
|--------|--------------|
| Larry Page | Google |
| Sergey Brin | Google |

### Disease-Outbreaks

| Disease | Time Period |
|---------|-------------|
| cholera | between 2010 and 2013 |

The Haiti **cholera** outbreak **between 2010 and 2013** was the worst epidemic of cholera in recent history.

### Person-Career

"This is not a victimless crime," said **Jim Kendall**, **president** of the Washington Association of Internet Service Providers.

| Person | Career |
|--------|--------|
| Jim Kendall | President |

### Man Made Disaster-Location

| Disaster | Location |
|----------|----------|
| fire | Booneville |

A **fire** destroyed a Cargill Meat Solutions beef processing plant in **Booneville**.

Other relations:
Person-Charge, Election-Winner, Natural Disaster-Location

Complex extraction systems:
CRFs, SVM kernels

Dense relations   Sparse relations

# Does Learning Ranking Models Help?



Person-Charge

Use **ranking model** based on F-measure of small set of queries

**Learn ranking model** on full document contents

- **Learning ranking models** leads to better document ranking
- **RSVM-IE** performs best at **early stages**
- **BAgg-IE** obtains high gains **later on**
- Objective function of **learning model** shapes document ranking

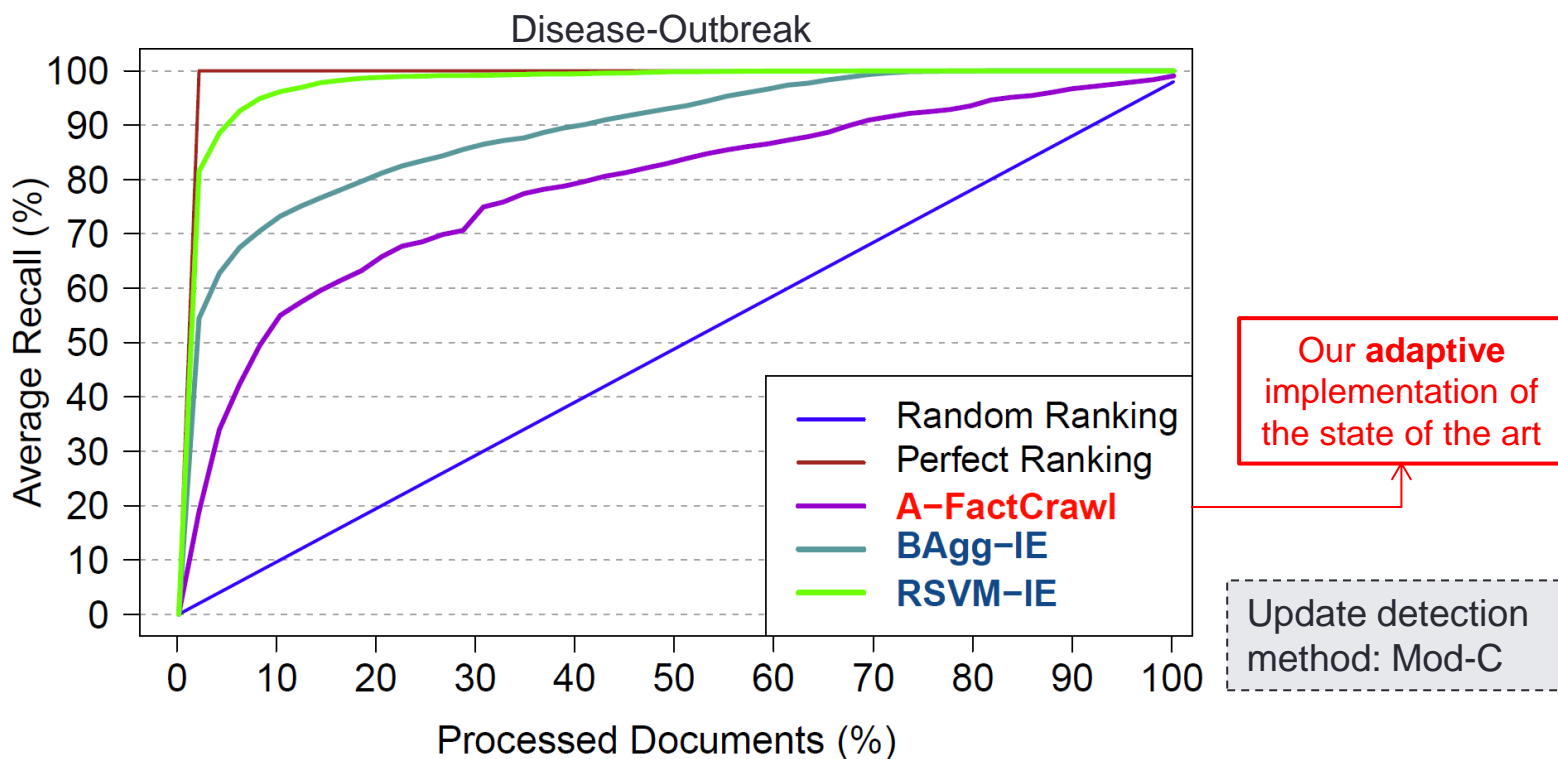Additional experiments in paper: analogous conclusions over all relations

# Does Update Detection Help?

Election-Winner



Update detection baselines:
**Wind-F**=Updates after processing 20,000 documents (2% of collection)
**Feat-S**=Update method based on Gaussian kernel *[Glazer, ICPR '12]*

Ranking Strategy: RSVM-IE

- **Feat-S** unable to evaluate over new features, crucial during adaptation
- **Top-*K*** and **Mod-C** improve the efficiency of the extraction process
- **Mod-C** leads to best execution using more efficient approach, with fewer models

Additional experiments in paper: analogous conclusions over all relations

# Putting Learning to Rank and Update Detection Together: Recall Analysis



Disease-Outbreak

Our **adaptive** implementation of the state of the art

Update detection method: Mod-C

- **Our techniques** bring **significant improvement** for sparse relations
- **RSVM-IE** **performs best**, as it prioritizes useful documents better, favoring adaptation

Additional experiments in paper: analogous conclusions over all relations

# Putting Learning to Rank and Update Detection Together: Extraction Time

**Person-Organization Affiliation**

Our **adaptive** implementation of the state of the art



- Cost of adapting in **A-FactCrawl** hurts efficiency of extraction process
- **Our techniques improve** efficiency of process even for inexpensive IE systems

Additional experiments in paper for **our techniques**:
- Analogous conclusions also for expensive IE systems and sparse relations
- **Scale linearly** in the size of the collection

# Document Ranking for Scalable Information Extraction: Summing Up

- Running IE system over **large** text collections is computationally **expensive**



Text Collection    IE system

<tornado, Florida>
...
<volcano, Chile>

- Proposed lightweight, adaptive approach and learning-based alternatives

  - **Online learning** algorithms with **in-training feature selection**: RSVM-IE, BAgg-IE

  

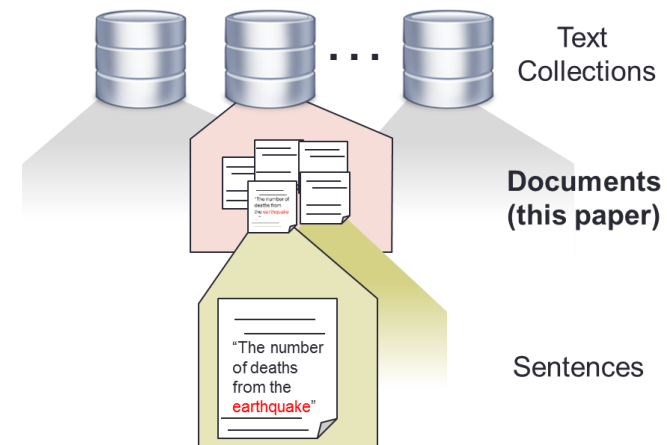  - **Update detection** based on feature changes: Mod-C, Top-*K*

  

- **RSVM-IE + Mod-C performs best**: Useful documents are better prioritized, enabling richer, more efficient ranking adaptation

# Future Work: Ranking at Different Granularities

- **Few collections** on the Web are relevant to an IE task

  Prioritize them based on number of useful documents

- **Few sentences** in a text document output tuples for an IE task

  Prioritize them based on usefulness and diversity



Text Collections

**Documents (this paper)**

"The number of deaths from the earthquake"

Sentences

# Future Work: Distributing the Execution of IE Systems

- Identify **optimal distributed execution strategy**

    E.g., by determining document placement in

    distributed file system



Document Collection

Map-Reduce Infrastructure

# But Before We Leave…

Try **REEL**, our toolkit to **easily develop and evaluate** IE systems

Open source and freely available at http://reel.cs.columbia.edu

Thanks!

# Information Extraction: Time Analysis

| Task | Time per sentence (ms) | | Toolkit or Algorithm | |
|---|---|---|---|---|
| Sentence splitting | 0.1 | | PTB | |
| Tokenization | 0.1 | | PTB | |
| Part-of-speech tagging | 7.4 | | ClearNLP | |
| Shallow parsing | 42 | | Search | |
| Dependency parsing | 25.6 | | ClearNLP | |
| Semantic role labeling | 8.4 | | ClearNLP | |
| Named Entity recognition (per entity) | 1.1 | | SENNA | |
| Relation extraction | 766 | 67 | Tree Kernel | OLLIE |
| Total | 850.7 | 151.7 | | |

# Experimental Settings: Data and Relations

- Dataset: **The New York Times** 1.8 million articles from 1987-2007

- Information Extraction Systems

Simple extraction systems:
HMMs, Text patterns

**Person-Organization**

**Google** co-founders **Larry Page** and **Sergey Brin** recently sat down with billionaire venture capitalist Vinod Khosla for a lengthy interview.

| Person | Organization |
|---|---|
| Larry Page | Google |
| Sergey Brin | Google |

**Disease-Outbreaks**

| Disease | Time Period |
|---|---|
| Cholera | between 2010 and 2013 |

The Haiti cholera outbreak **between 2010 and 2013** was the worst epidemic of **cholera** in recent history.

**Person-Career**

"This is not a victimless crime," said **Jim Kendall**, **president** of the Washington Association of Internet Service Providers.

| Person | Career |
|---|---|
| Jim Kendall | President |

**Man Made Disaster-Location**

| Disaster | Location |
|---|---|
| fire | Booneville |

A **fire** destroyed a Cargill Meat Solutions beef processing plant in **Booneville**.

**Person-Charge**

| Person | Charge |
|---|---|
| Ibrahim Muktar Said | Connection with bombing |

**Ibrahim Muktar Said** was charged Sunday night in **connection with** the failed Hackney bus **bombing**.

**Election-Winner**

| Person | Election |
|---|---|
| Boris Johnson | London mayoral election |

**Boris Johnson** defeated Ken Livingstone in the **London mayoral election**.

**Natural Disaster-Location**

| Disaster | Location |
|---|---|
| tornado | Florida |

A **tornado** swept the coast of **Florida** on Wednesday.

Dense relations

Complex extraction systems:
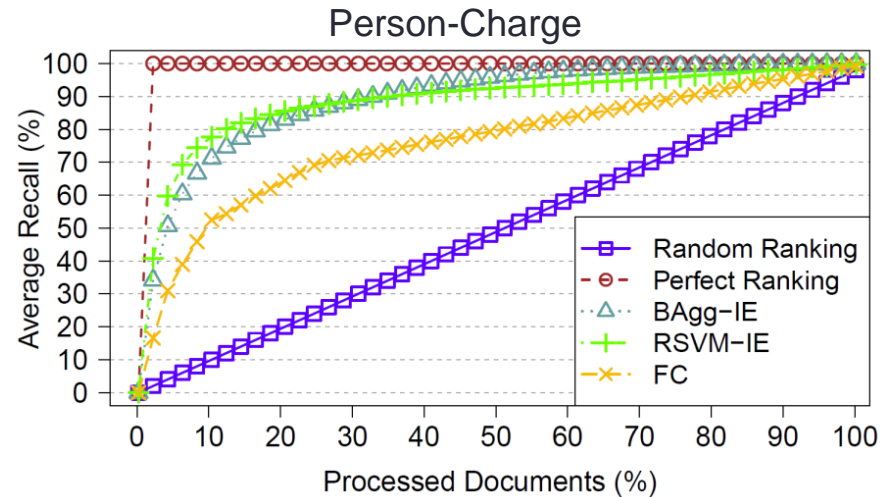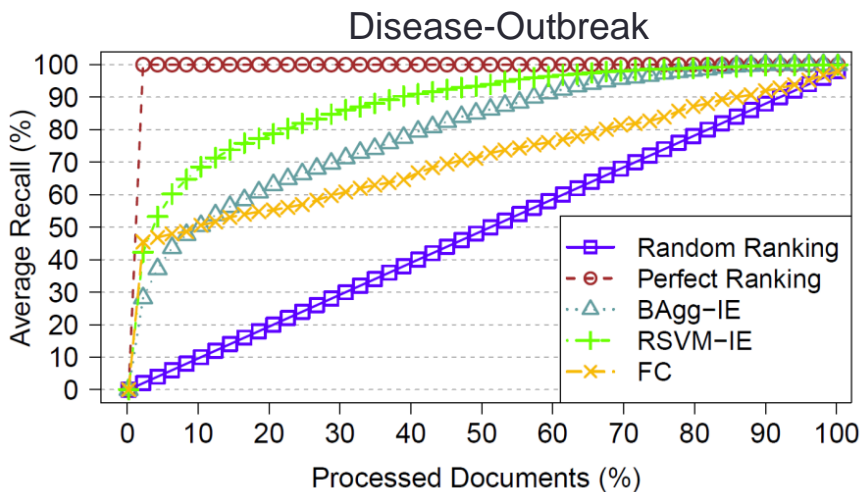CRFs, SVM Kernels

Sparse relations

# Experimental Settings: Extractors

- Person-Organization Affiliation:
  - Entities: HMM and text patterns
  - Relation: SVM classifier

- Disease-Outbreak:
  - Entities: Dictionaries and manually crafted regular expressions
  - Relation: Distance between entities

- Others:
  - Entities: Stanford NLP (Person and Location), MEMM (Natural Disasters), and CRF (others)
  - Relation: Subsequences Kernel *[Bunescu and Mooney, NIPS '05]*
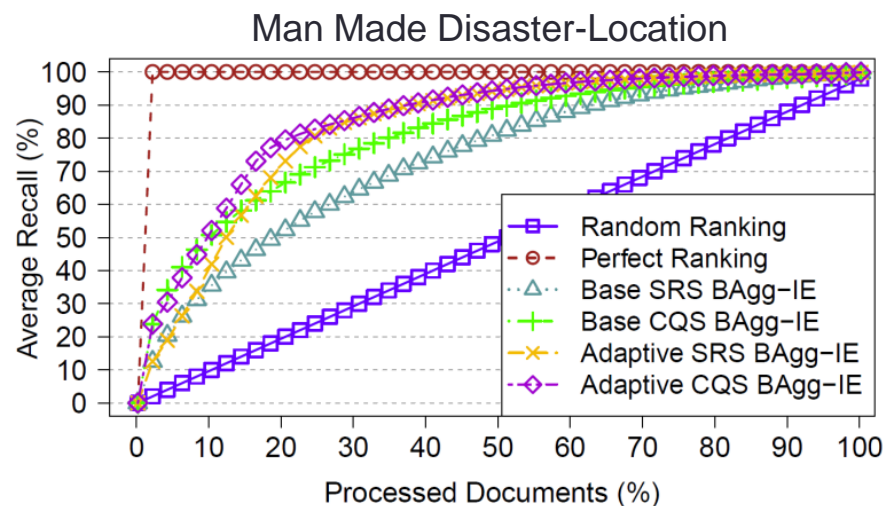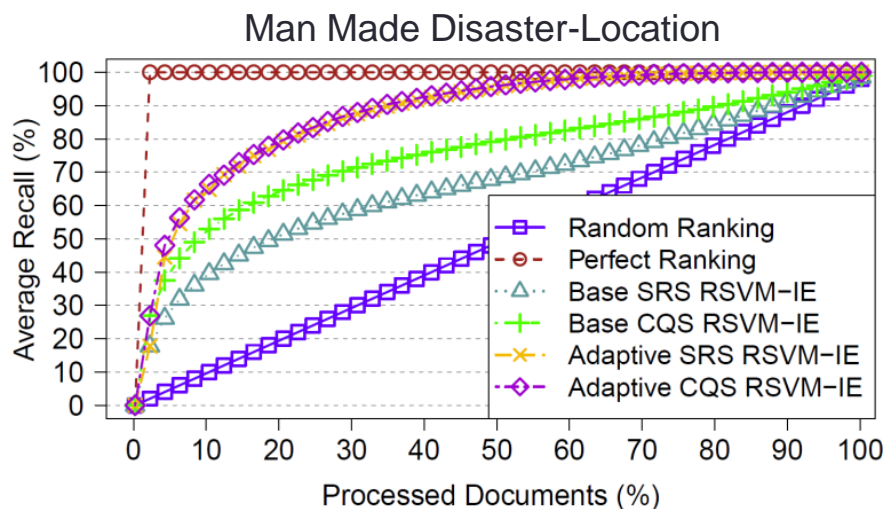
# Experimental Settings: Details

- Document Sampling Strategies:

  - Simple Random Sampling (**SRS**): Documents are collected randomly from fully-accessible collection

  - Cyclic Querying Sampling (**CQS**): Queries learned from external collection and issued in a round-robin fashion

- Update Detection:

  - Feature Shifting (**Feat-S**): Gaussian kernel for one-class classification
    - Triggers an update for high geometrical difference
      [A. Glazer, "Feature Shift Detection." *ICPR* '12]

  - Fixed Window (**Wind-F**): Triggers after processing N documents
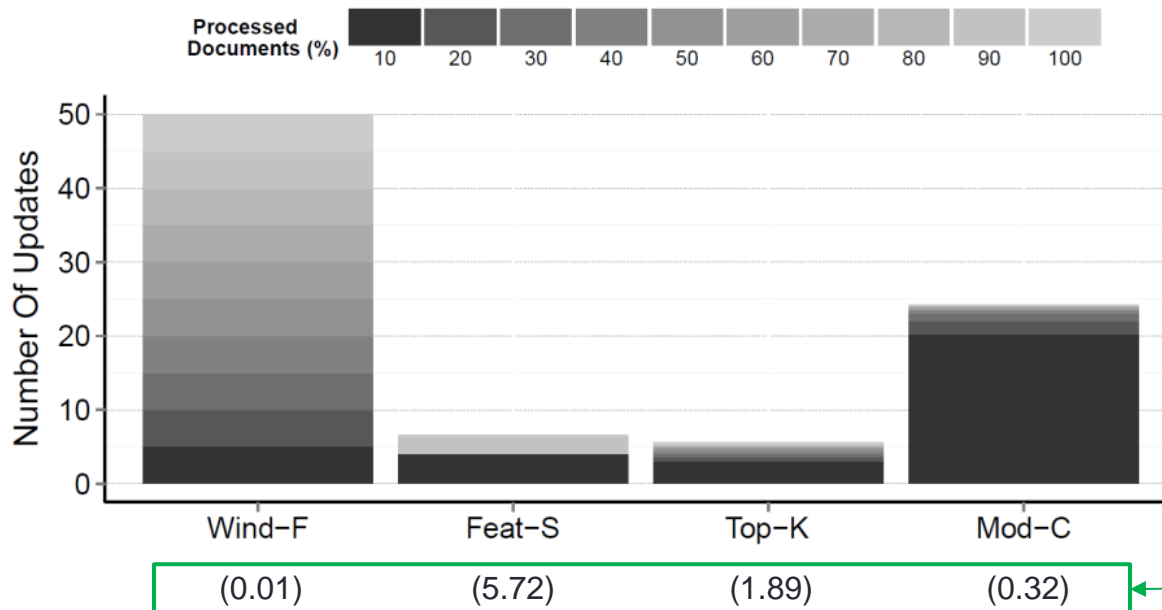
# Ranking Models vs. FactCrawl



- Using **full document contents** leads to better document ranking
- RSVM-IE performs best at **early stages**
- BAgg-IE obtains high gains **later on**
- **Objective function shapes the document ranking**

# Impact of Document Sampling



Man Made Disaster-Location

Man Made Disaster-Location

- CQS improves recall at early stages
- CQS obtains higher average precision and AUC
- Targeted sampling **improves the efficiency of** the extraction process

# Update Detection: Time and Distribution of Updates



Update detection baselines:
**Wind-F**=Updates after processing 20,000 documents (2% of collection)
**Feat-S**=Update method based on Gaussian kernel *[Glazer, ICPR '12]*
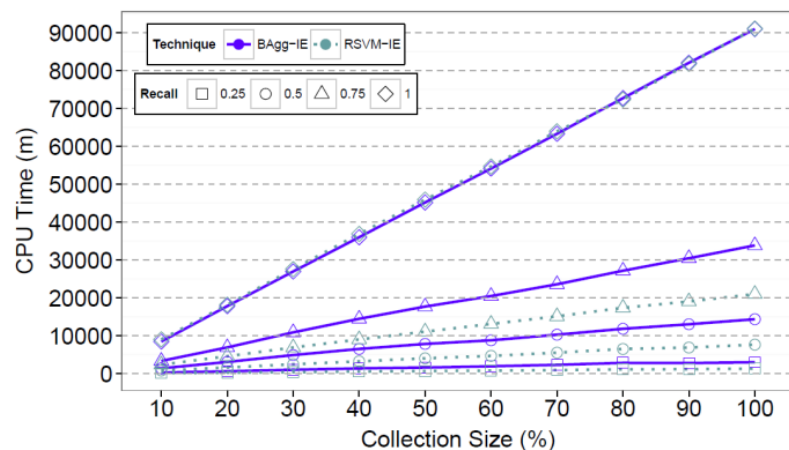
Average CPU time per document (ms)

- Wind-F is the **most efficient** but ignores document contents
- Feat-S performs fewer updates but is affected by kernel cost
- Top-*K* performs the **fewest** updates, relatively efficiently
- Mod-C exhibits best **number of updates-time balance**
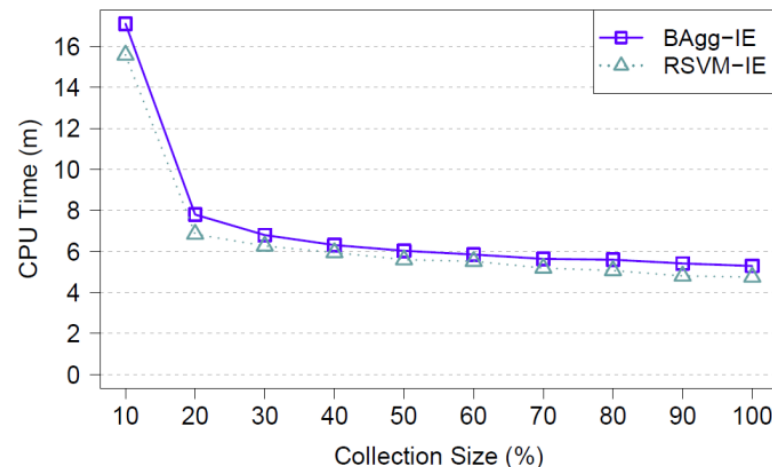
# Scalability Analysis: Running Time

**Target recall**

Natural Disaster-Location

**Fixed set of useful documents**

Person-Organization Affiliation



- Our approach:
  - **Scales linearly** to collection size
  - **Improves** with the more information we find in larger collections
  - **Is a substantial step towards scalable information extraction**