# Optimizing Scoring Functions and Indexes for Proximity Search in Type-annotated Corpora

Soumen Chakrabarti
▶ Kriti Puniyani
Sujatha Das

IIT Bombay

---

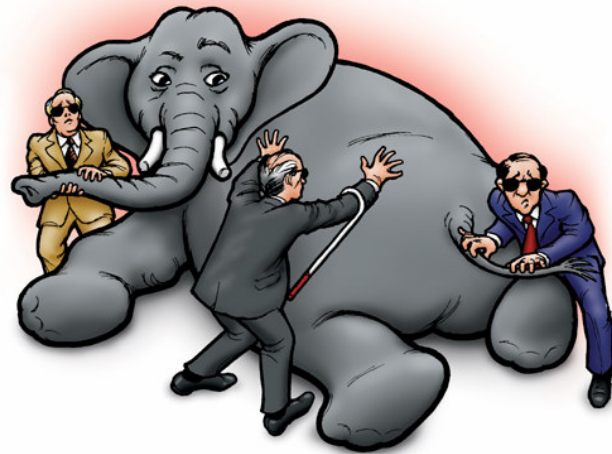(In fewer words)

# Ranking and Indexing for Semantic Search

Soumen Chakrabarti
▶ Kriti Puniyani
Sujatha Das

IIT Bombay

# Working notion of semantic search

- Exploiting in conjunction
  - "Strings with meaning" – entities and relations
  - "Uninterpreted strings" – as in IR
- This paper
  - Only "is-a" relation
  - Token match
  - Token proximity
- Can approximate many info needs



# Type-annotated corpus and query e.g.

Name a **physicist** who searched for intelligent life in the cosmos
→ type=physicist NEAR "cosmos"…

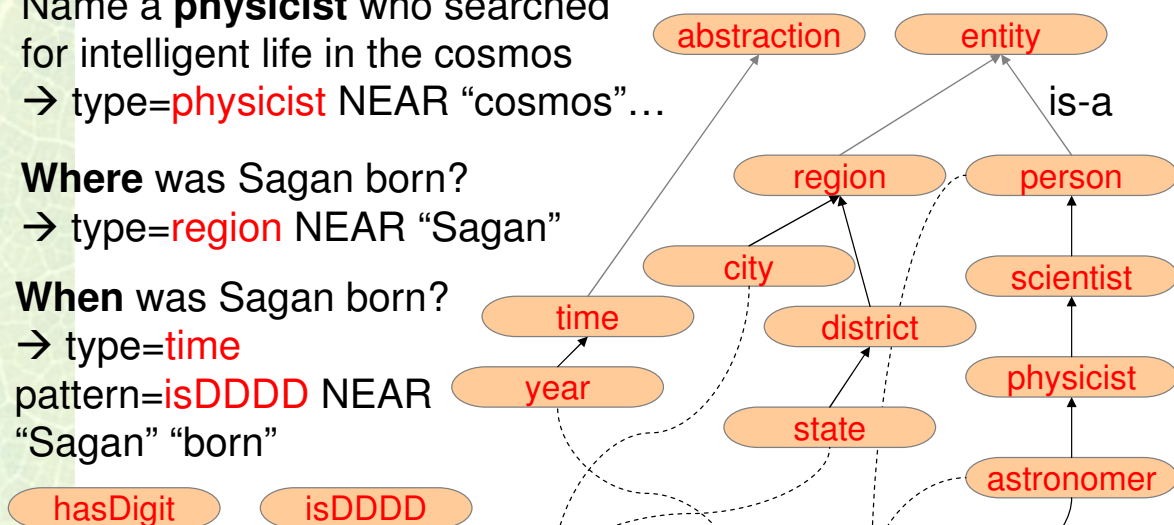**Where** was Sagan born?
→ type=region NEAR "Sagan"

**When** was Sagan born?
→ type=time
pattern=isDDDD NEAR
"Sagan" "born"



is-a

Born in New York in 1934 , Sagan was
a noted astronomer whose lifelong passion
was searching for intelligent life in the cosmos.

4

# The query class we address

- Find a token span *w* (in context) such that
  - *w* is a mention of entity *e*
    - "Carl Sagan" or "Sagan" is a mention of the concept of that specific physicist
  - *e* is an instance of **atype** *a* given in the query
    - Which *a*=physicist …
  - *w* is "NEAR" a set of **selector** strings
    - "searched", "intelligent", "life", "cosmos"
- All uncertain/imprecise; we focus on #3
- Yet surprisingly powerful: correct answer within top 3—4 *w*'s for TREC QA benchmark

# Contribution 1: What is "NEAR"?

- XQuery and XPath full text support
  - (distance at most|window) 10 words [ordered] – hard proximity clause, not learnt
  - ftcontains … with thesaurus at … relationship "narrower terms" at most $\ell$ levels
- No implementation combining "narrower terms" and "soft" proximity ranking
- Search engines favor proximity in proprietary ways
- A learning framework for proximity

# Contribution 2: Indexing annotations

- type=person NEAR theory relativity → type in {physicist, politician, cricketer,…} NEAR theory relativity
  - Large fanout at query time, impractical
- Complex annotation indexes tend to be large
  - Binding Engine (WWW 2005): 10x index size blowup with only a handful of entity types
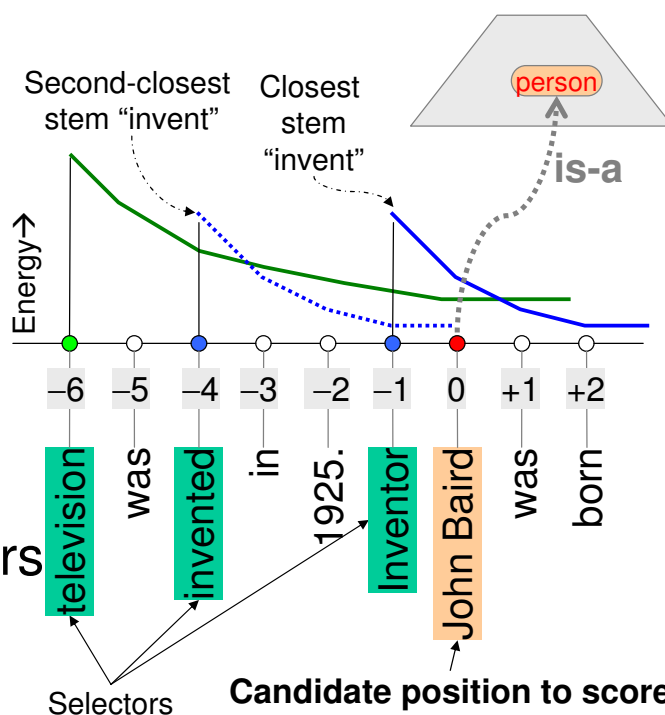  - Our target: 18000 atypes today, more later
- Workload-driven index and query optimization
  - Exploit skew in query atype workload

# Part-1: Learning to score token spans

- type=person NEAR "television" "invent*"
- Rarity of selectors
- Distance from candidate position to selectors
- Many occurrences of one selector
  - Closest is good
- Combining scores from many selectors
  - Sum is good



Second-closest stem "invent"

Closest stem "invent"

person

is-a

Energy→

−6 −5 −4 −3 −2 −1 0 +1 +2

television  was  invented  in  1925.  Inventor  John Baird  was  born

Selectors

**Candidate position to score**

# Learning the shape of the decay function

- For simplicity assume left-right symmetry
- Parameters $(\beta_1,\ldots,\beta_W)$, $W$=max gap window
- Candidate position characterized by a feature vector $f = (f[1],\ldots,f[W])$
  - If there is a matched selector $s$ at distance $j$ and
  - This is the closest occurrence of $s$
  - Then set $f[j]$ to $energy(s)$, … else 0
- Score of candidate position is $\beta \cdot f$
- If we like candidate $u$ less than $v$ ("$u \prec v$")
  - We want $\beta \cdot f_u \leq \beta \cdot f_v$
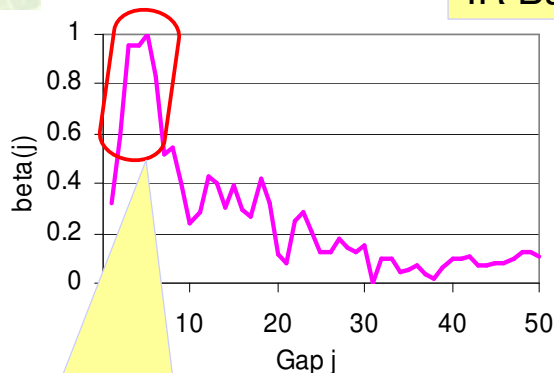  - Assess a penalty proportional to $\exp(\beta \cdot f_u - \beta \cdot f_v)$

9

---

# Learning decay function—results

$$\min_{\beta} \sum_{j=1}^{W} (\beta_j - \beta_{j+1})^2 + C \sum_{u \prec v} \exp(\beta \cdot f_u - \beta \cdot f_v)$$

Discourage adjacent βs from differing a lot

Penalize violations of preference order

IR Baseline



beta (j) vs Gap j

Roughly unimodal around gap = 4 and 5

| Train | Test | MRR |
|-------|------|-----|
| IR | 2000 | 0.16 |
| 2001 | 2000 | **0.29** |

TREC year

Mean reciprocal rank: Average over questions, reciprocal of the first rank where an answer token was found (large good)

10

# Part-2: Workload-driven indexing

- Type hierarchies are large and deep
  - 18000 internal and 80000 leaf types in WordNet
- Runtime atype expansion time-intensive
  - Even WordNet knows 650 scientists, 860 cities…
- Index each token as all generalizations
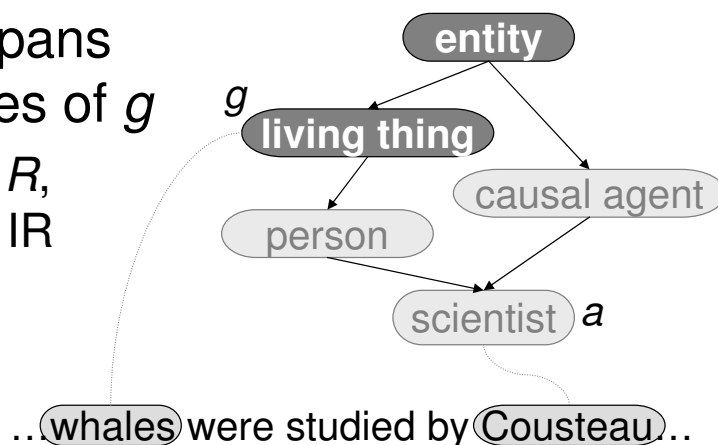  - Sagan → physicist, scientist, person, living thing
  - Large index space bloat
- 💡Index a subset of atypes

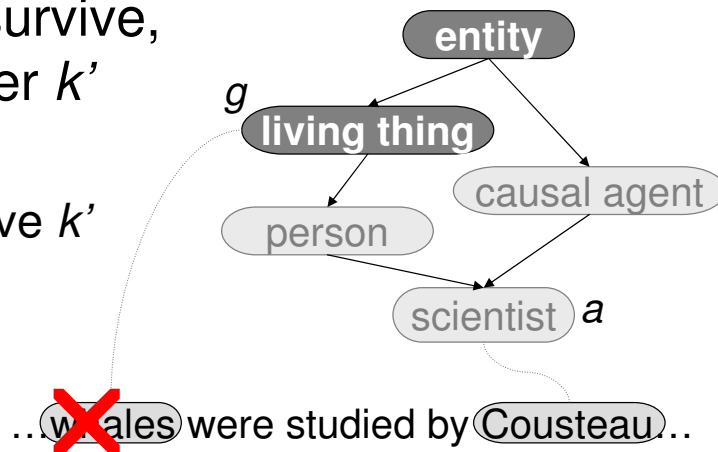| Corpus/Index | Gbytes |
|---|---|
| Original corpus | 5.72 |
| Gzipped corpus | 1.33 |
| Stem index | 0.91 |
| Full type index | 4.30 |

# Pre-generalize (and post-filter)

- Full set of "atypes" (answer types) is *A*
- Index only a "registered" subset *R* of *A*
- Say query has atype *a*; want *k* answers
- Find *a*'s "best" generalization *g*∈ *R*
- Get best *k'* >*k* spans that are instances of *g*
  - Given index on *R*, this is standard IR (see paper)



*g*

entity

living thing

causal agent

person

scientist *a*

..whales were studied by Cousteau..

# (Pre-generalize and) post-filter

- Fetch each high-scoring span $w$
- Check if $w$ is-a $a$
  - Fast compact "forward index" (doc,offset)→token
  - Fast small "reachability index", common in XML
- If fewer than $k$ survive, restart with larger $k'$
  - Expensive
  - Pick conservative $k'$

entity

$g$

living thing

causal agent

person

scientist $a$

✖

...whales were studied by Cousteau...

13

---

# Estimates needed by optimizer

- If we index token ancestors in $R$ as against ancestors in all of $A$, how much index space will we save?
  - Cannot afford to try out and see for many $R$s

- If query atype $a$ is not found in $R$ and we must generalize to $g$, what will be the bloat factor in query processing time?
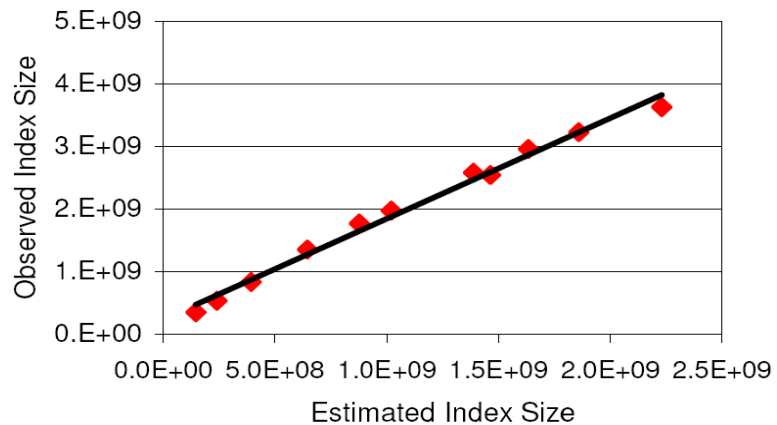  - Need to average over a representative workload

14

# Index space estimate given R

- Each token occurrence leads to one posting entry
- Assume index compression is a constant factor
- Then total estimated index size is proportional to

$$\sum_{r \in R} corpusCount(r)$$

Number of tokens in corpus that connect up to $r$

- Surprisingly accurate!

(Chart: vertical axis "Observed Index Size" with values 0.E+00, 1.E+09, 2.E+09, 3.E+09, 4.E+09, 5.E+09; horizontal axis "Estimated Index Size" with values 0.0E+00, 5.0E+08, 1.0E+09, 1.5E+09, 2.0E+09, 2.5E+09)

---

# Processing time bloat for one query

- If $R=A$, query takes time approximated by

$$t_{scan}\, corpusCount(a)$$

Time to score one candidate position while scanning postings

Number of occurrences of descendants of type $a$

- If $a$ cannot be found in $R$, the price paid for generalization to $g$ consists of
  - Scanning more posting entries: $t_{scan}\, corpusCount(g)$
  - Post-filtering $k'$ responses: $k'\, t_{filter}$
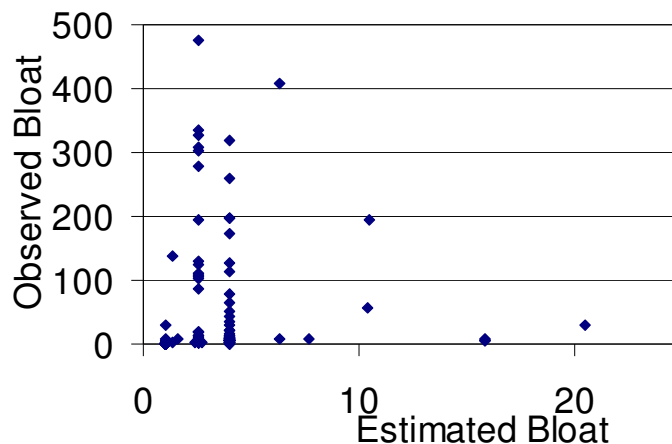- Therefore, overall bloat factor is

$$queryBloat(a,R) = \frac{t_{scan}\, corpusCount(g) + k'\, t_{filter}}{t_{scan}\, corpusCount(a)}$$

Time to check if answer is instance of $a$ as well

16

# Query time bloat—results



- Observed bloat fit not as good as index space estimate

- While observed::estimated ratio for one query is noisy, average over many queries is much better

---

# Expected bloat over many queries

Prob of new query having atype $a$

$$\sum_{a \in A} queryProb(a)\ queryBloat(a, R)$$

Already estimated

- Maximum likelihood estimate

$$queryProb_{\text{Train}}(a) = \frac{queryCount_{\text{Train}}(a)}{\sum_{a' \in A} queryCount_{\text{Train}}(a')}$$

- Many $a$'s get zero training probability
  → Optimizer does not register $g$ close to $a$
- Low-prob atypes appear in test → huge bloat
- Collectively matter a lot (heavy-tailed distrib)
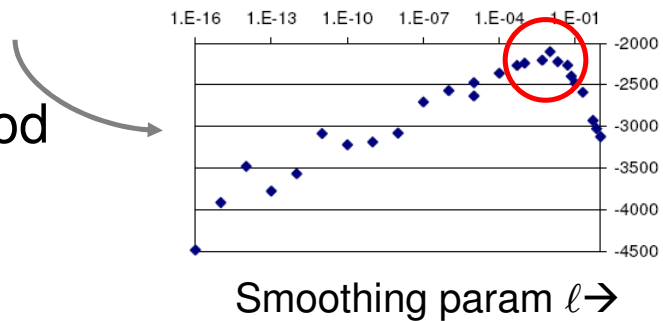
# Smoothing low-probability atypes

- Lidstone smoothing:

$$queryProb_{Train}(a) = \frac{queryCount_{Train}(a) + \ell}{\sum_{a' \in A}\left(queryCount_{Train}(a') + \ell\right)}$$

- Smoothing param $\ell$ fit by maximizing log-likelihood of held-out data:

$$\sum_{a \in HeldOut} queryCount_{HeldOut}(a)\log(queryProb_{Train}(a))$$

- Clear range of good fits for $\ell$



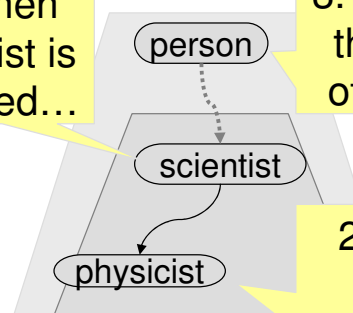Smoothing param $\ell\rightarrow$

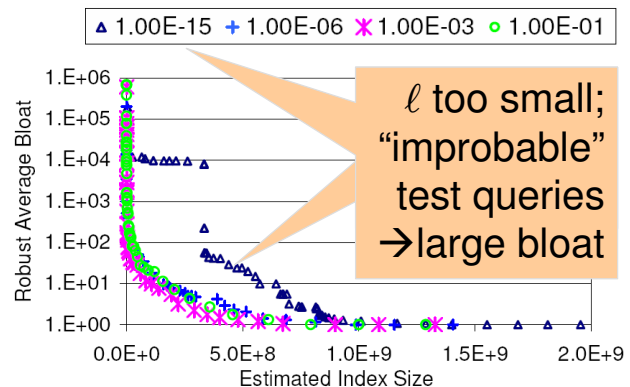# The *R* selection algorithm

- $R \leftarrow$ roots of $A$
- Greedily add the "most profitable" atype $a^*$
- Profit = ratio of
  - reduction in bloat of $a^*$ and its descendants to
  - increase in index space
- Downward and upward traversals and updates
- Gives a tradeoff between index space and query bloat
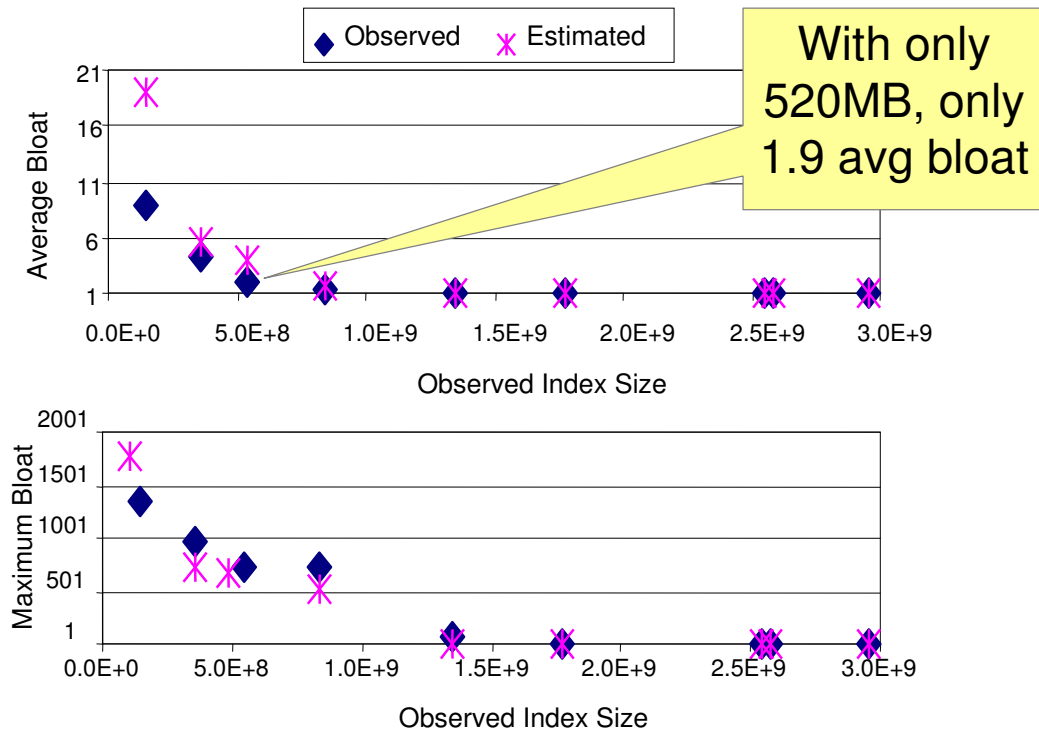


1. When scientist is included…

3. reducing the profit of person

2. bloat of physicist goes down

$\ell$ too small; "improbable" test queries →large bloat

# Optimized space-time tradeoff



With only 520MB, only 1.9 avg bloat

# Optimized index sizes

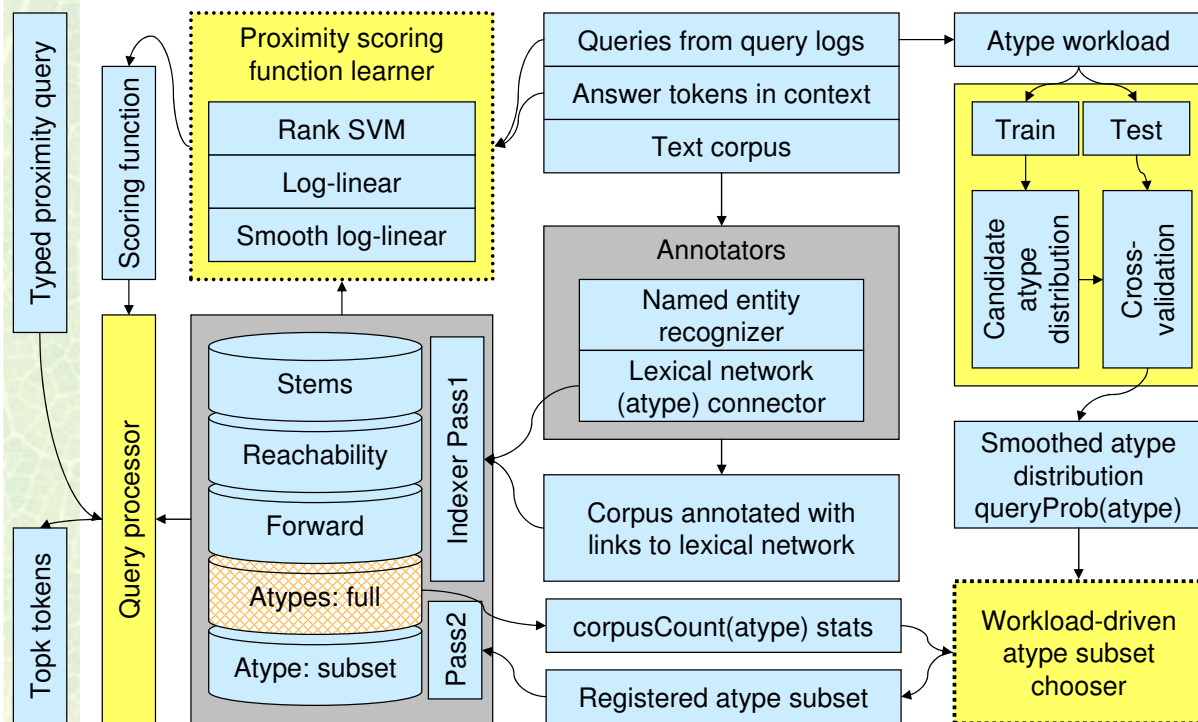| Corpus/Index | Gbytes |
|---|---|
| Original corpus | 5.72 |
| Gzipped corpus | 1.33 |
| Stem index | 0.91 |
| Full type index | 4.30 |
| Reachability index | 0.01 |
| Forward index | 1.16 |
| Atype subset index | 0.52 |

# Summary

- **Working prototype around Lucene and UIMA**
  - Annotators attach tokens to type taxonomy
  - Query atype workload help compact index
  - Ranking function learnt from preference data
  - NL queries translated into atype+selectors
- **Ongoing work**
  - Indexing and searching relations other than is-a
  - More general notions of graph proximity
- **Email soumen@cse.iitb.ac.in for code access**

# The big picture



Email soumen@cse.iitb.ac.in for code access