# Introduction to
# Information Retrieval

**Systems issues**

# Background

- Score computation is a large (10s of %) fraction of the CPU work on a query
    - Generally, we have a tight budget on latency (say, 250ms)
    - CPU provisioning doesn't permit exhaustively scoring every document on every query
- Today we'll look at ways of cutting CPU usage for scoring, without compromising the quality of results (much)
- Basic idea: avoid scoring docs that won't make it into the top $K$

# Safe vs non-safe ranking

- The terminology "safe ranking" is used for methods that guarantee that the $K$ docs returned are the $K$ absolute highest scoring documents

- Is it ok to be non-safe?

# Ranking function is only a proxy

- User has a task and a query formulation
- Ranking function matches docs to query
- Thus the ranking function is anyway a proxy for user happiness
- If we get a list of $K$ docs "close" to the top $K$ by the ranking function measure, should be ok

# Recap: Queries as vectors

- **Key idea 1:** Do the same for queries: represent them as vectors in the space

- **Key idea 2:** Rank documents according to their proximity to the query in this space

- proximity = similarity of vectors, measured by cosine similarity

# Efficient cosine ranking

- Find the $K$ docs in the collection "nearest" to the query $\Rightarrow K$ largest query-doc cosines.
- Efficient ranking:
  - Computing a single cosine efficiently.
  - Choosing the $K$ largest cosine values efficiently.
    - Can we do this without computing all $N$ cosines?

# Computing the *K* largest cosines: selection vs. sorting

- Typically we want to retrieve the top *K* docs (in the cosine ranking for the query)
  - not to totally order all docs in the collection
- Can we pick off docs with *K* highest cosines?
- Let *J* = number of docs with nonzero cosines
  - We seek the *K* best of these *J*

# Bottlenecks

- Primary computational bottleneck in scoring: <u>cosine computation</u>

- Can we avoid all this computation?

- Yes, but may sometimes get it wrong
  - a doc *not* in the top *K* may creep into the list of *K* output docs
  - As noted earlier, this may not be a bad thing

# SPEEDING COSINE COMPUTATION BY PRUNING

# Generic approach

- Find a set *A* of *contenders*, with *K* < |*A*| << *N*
  - *A* does not necessarily contain the top *K,* but has many docs from among the top *K*
  - Return the top *K* docs in *A*
- Think of *A* as <u>pruning</u> non-contenders
- The same approach is also used for other (non-cosine) scoring functions

# Index elimination

- Basic cosine computation algorithm only considers docs containing at least one query term

- Take this further:

  - Only consider high-idf query terms

  - Only consider docs containing many query terms

# High-idf query terms only

- For a query such as *catcher in the rye*
- Only accumulate scores from *catcher* and *rye*
- Intuition: **in** and **the** contribute little to the scores and so <u>don't alter rank-ordering much</u>
- Benefit:
  - Postings of low-idf terms have many docs → these (many) docs get eliminated from set *A* of contenders

# Docs containing many query terms

- Any doc with at least one query term is a candidate for the top *K* output list

- For multi-term queries, only compute scores for docs containing several of the query terms
  - Say, at least 3 out of 4
  - Imposes a "soft conjunction" on queries seen on web search engines (early Google)

- Easy to implement in postings traversal

# 3 of 4 query terms

| *Antony* | → | 3 | 4 | 8 | 16 | 32 | 64 | 128 | |
| *Brutus* | → | 2 | 4 | 8 | 16 | 32 | 64 | 128 | |
| *Caesar* | → | 1 | 2 | 3 | 5 | 8 | 13 | 21 | 34 |
| *Calpurnia* | → | 13 | 16 | 32 | | | | | |

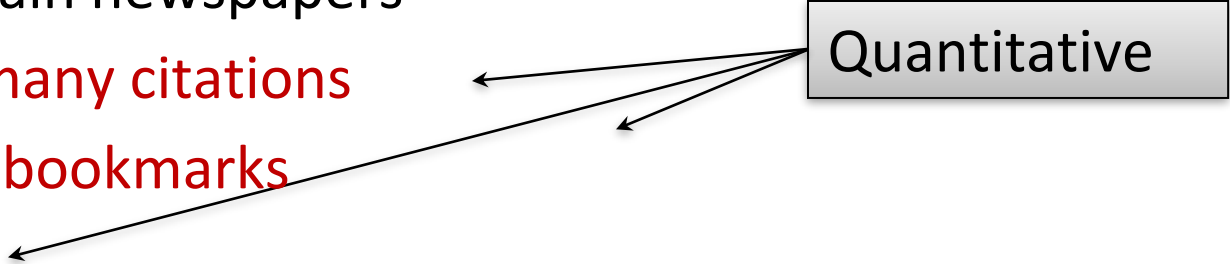Scores only computed for docs 8, 16 and 32.

# Champion lists

- Precompute for each dictionary term $t$, the $r$ docs of highest weight in $t$'s postings
    - Call this the <u>champion list</u> for $t$
    - (aka <u>fancy list</u> or <u>top docs</u> for $t$)
- Note that $r$ has to be chosen at index build time
    - Thus, it's possible that $r < K$
- *Highest tf among docs*
- At query time, only compute scores for docs in the champion list of some query term
    - Pick the $K$ top-scoring docs from amongst these

# QUERY-INDEPENDENT DOCUMENT SCORES

# Static quality scores

- We want top-ranking documents to be both *relevant* and *authoritative*

- *Relevance* is being modeled by cosine scores

- *Authority* is typically a query-independent property of a document

- Examples of authority signals
  - Wikipedia among websites
  - Articles in certain newspapers
  - A paper with many citations
  - Many likes, or bookmarks
  - Pagerank

Quantitative

# Modeling authority

- Assign to each document *d* a *query-independent* <u>quality score</u> in [0,1]
  - Denote this by *g(d)*
- <span style="color:red">Thus, a quantity like the number of citations is scaled into [0,1]</span>

# Net score

- Consider a simple total score combining cosine relevance and authority
- net-score(*q,d*) = *g(d)* + cosine(*q,d*)
  - Can use some other linear combination
  - Indeed, any function of the two "signals" of user happiness
- Now we seek the top *K* docs by <u>net score</u>

# Top $K$ by net score – fast methods

- First idea: Order all postings by $g(d)$
- Key: this is a common ordering for all postings
- Thus, can concurrently traverse query terms' postings for
    - Postings intersection
    - Cosine score computation

# Why order postings by *g(d)?*

- Under *g(d)*-ordering, top-scoring docs likely to appear early in postings traversal

- In time-bound applications (say, we have to return whatever search results we can in 50 ms), this allows us to stop postings traversal early

  - Short method of computing scores for all docs in postings

# Champion lists in *g(d)*-ordering

- Can combine champion lists with *g(d)*-ordering

- <span style="color:red">Maintain for each term a champion list of the *r* docs with highest *g(d)* + tf-idf$_{td}$</span>

- Seek top-*K* results from only the docs in these champion lists
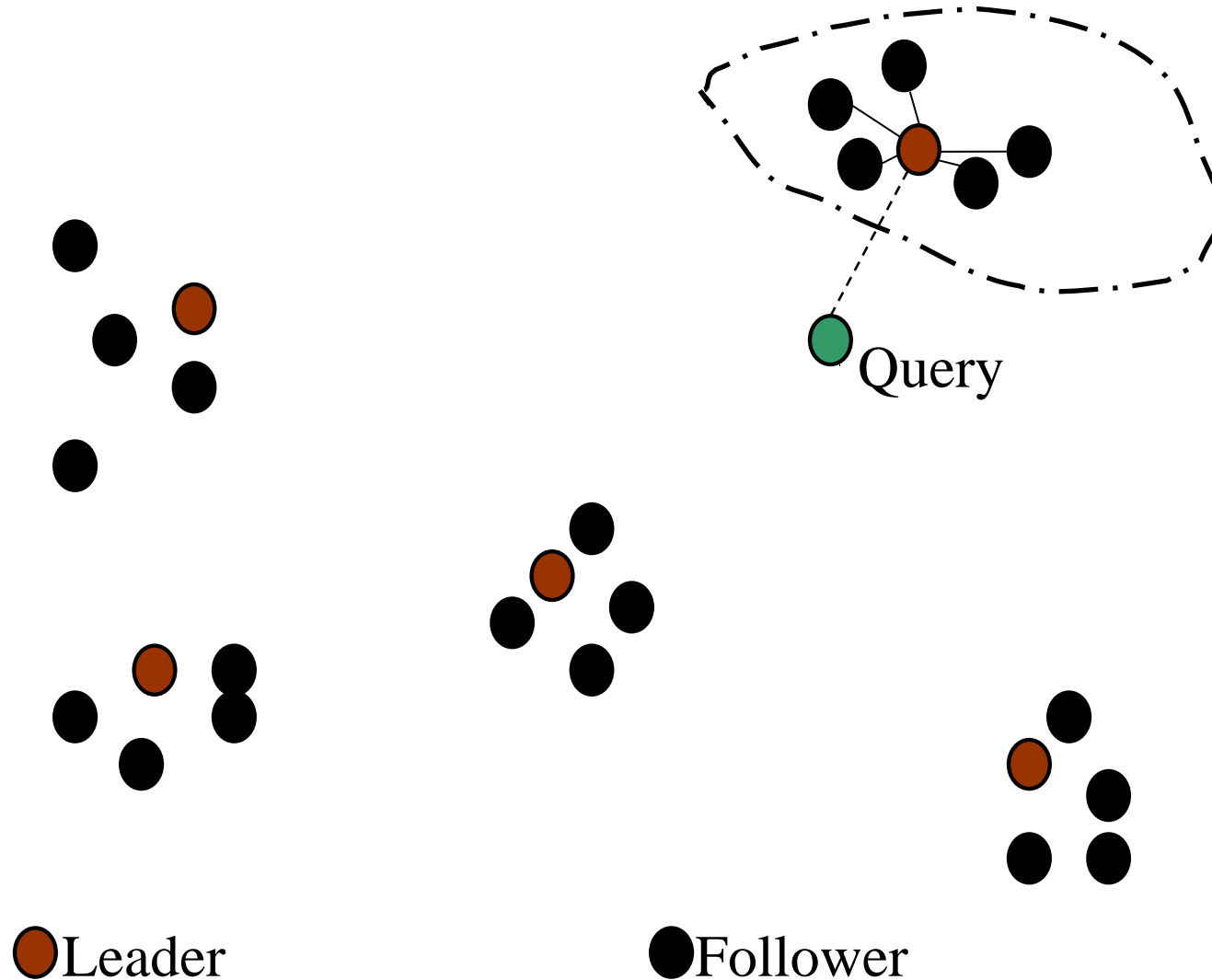
# CLUSTER PRUNING

# Cluster pruning: preprocessing

- Pick √N *docs* at random: call these *leaders*
- For every other doc, pre-compute nearest leader
  - Docs attached to a leader: its *followers;*
  - <u>Likely</u>: each leader has ~ $\sqrt{N}$ followers.

# Cluster pruning: query processing

- Process a query as follows:
  - Given query *Q*, find its nearest *leader L.*
  - Seek *K* nearest docs from among *L*'s followers.

# Visualization



Query

Leader    Follower

# Why use random sampling

- Fast

- Leaders likely to reflect data distribution (try not to be biased)

# Impact-ordered postings

- We only want to compute scores for docs for which $tf_{t,d}$ is high enough

- We sort each postings list by $tf_{t,d}$

- Now: not all postings in a common order (as per doc id)

- How do we compute scores in order to pick off top *K?*

  - Two ideas follow

# 1. Early termination

- When traversing *t's* postings, stop early after either
  - a fixed number of *r* docs
  - $tf_{t,d}$ drops below some threshold
- <span style="color:darkred">Take the union of the resulting sets of docs</span>
  - <span style="color:darkred">One from the postings of each query term</span>
- Compute only the scores for docs in this union

# 2. idf-ordered terms

- When considering the postings of query terms
- Look at them in order of decreasing idf
  - High idf terms likely to contribute most to score
- As we update score contribution from each query term
  - Stop if doc scores relatively unchanged
- Can apply to cosine or some other net scores

# TIERED INDEXES

# High and low lists

- For each term, we maintain two postings lists called *high* and *low*
    - Think of *high* as the champion list
- When traversing postings on a query, only traverse *high* lists first
    - If we get more than *K* docs, select the top *K* and stop
    - Else proceed to get docs from the *low* lists
- Can be used even for simple cosine scores, without global quality *g(d)*
- A means for segmenting index into two <u>tiers</u>

# Tiered indexes

- Break postings up into a hierarchy of lists
  - Most important
  - …
  - Least important
- Can be done by *g(d)* or another measure
- Inverted index thus broken up into <u>tiers</u> of decreasing importance
- At query time use top tier unless it fails to yield *K* docs
  - If so drop to lower tiers

# Example tiered index