

# Web Search: Techniques, algorithms and Applications

## Basic Techniques for Web Search

German Rigau <[german.rigau@ehu.es](mailto:german.rigau@ehu.es)>

[Based on slides by Eneko Agirre ...  
and Christopher Manning and Prabhakar Raghavan]



# Basic Techniques for Web Search

- Review of applications
- Basic Techniques in detail:
  - Boolean search
  - Vocabularies, dictionaries, index
  - **Scoring, evaluation, complete system**
  - Web search
- Semantic search

# Complete system (Chap. 7)

- Putting together a complete search system
  - Will require learning about a number of miscellaneous topics and heuristics

# 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 diggs, Y!buzzes or del.icio.us marks
  - (Pagerank)



Quantitative

# Modeling authority

- Assign to each document a *query-independent* quality score in  $[0,1]$  to each document  $d$ 
  - Denote this by  $g(d)$
- Thus, a quantity like the number of citations is scaled into  $[0,1]$

# Net score

- Consider a simple total score combining cosine relevance and authority
- $\text{net-score}(q,d) = g(d) + \text{cosine}(q,d)$ 
  - Can use some other linear combination than an equal weighting
  - Indeed, any function of the two “signals” of user happiness
- Now we seek the top  $K$  docs by net score

# 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 of computing scores for all docs in postings



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

- Can combine champion lists with  $g(d)$ -ordering
- Maintain for each term a champion list of the  $r$  docs with highest  $g(d) + \text{tf-idf}_{td}$
- Seek top- $K$  results from only the docs in these champion lists

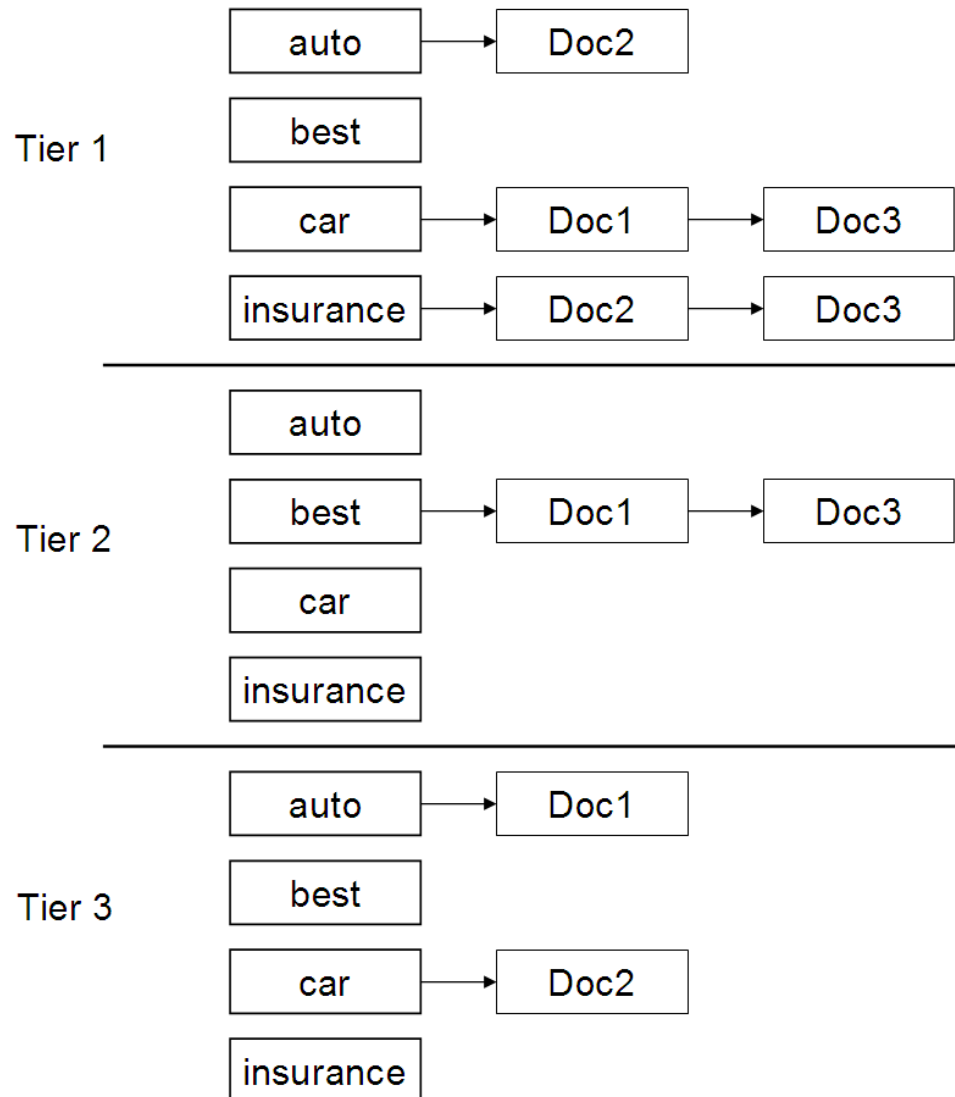
# High and low lists

- 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 tiers

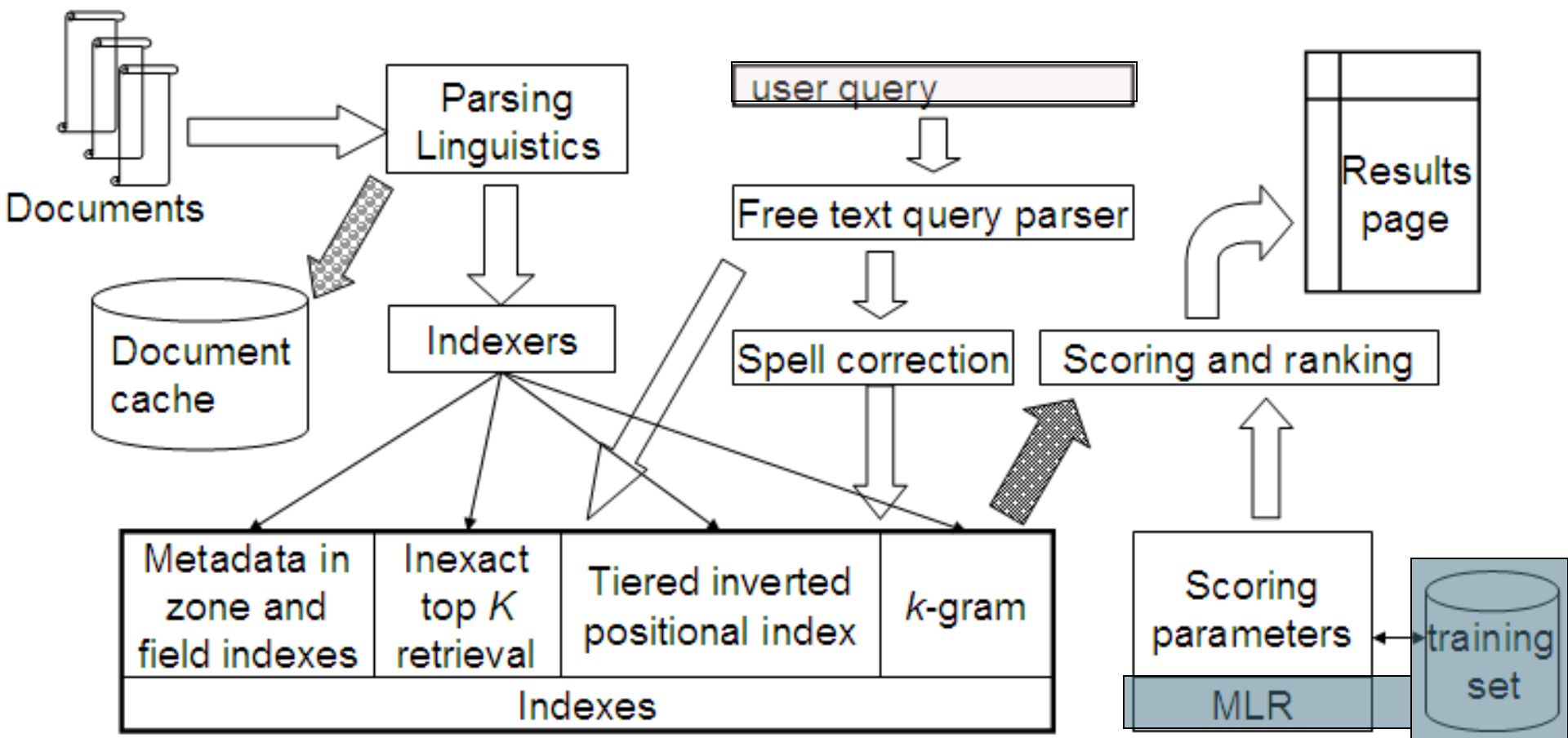
# Tiered indexes

- 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 tiers 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



# Putting it all together



# Results presentation

■ ...

# Result Summaries

- Having ranked the documents matching a query, we wish to present a results list
- Most commonly, a list of the document titles plus a short summary, aka “10 blue links”

## [John McCain](#)

**John McCain** 2008 - The Official Website of **John McCain's** 2008 Campaign for President ... African American Coalition; Americans of Faith; American Indians for **McCain**; Americans with ...  
[www.johnmccain.com](http://www.johnmccain.com) · [Cached page](#)

## [JohnMcCain.com - McCain-Palin 2008](#)

**John McCain** 2008 - The Official Website of **John McCain's** 2008 Campaign for President ... African American Coalition; Americans of Faith; American Indians for **McCain**; Americans with ...  
[www.johnmccain.com/Informing/Issues](http://www.johnmccain.com/Informing/Issues) · [Cached page](#)

## [John McCain News- msnbc.com](#)

Complete political coverage of **John McCain**. ... Republican leaders said Saturday that they were worried that Sen. **John McCain** was heading for defeat unless he brought stability to ...  
[www.msnbc.msn.com/id/16438320](http://www.msnbc.msn.com/id/16438320) · [Cached page](#)

## [John McCain | Facebook](#)

Welcome to the official Facebook Page of **John McCain**. Get exclusive content and interact with **John McCain** right from Facebook. Join Facebook to create your own Page or to start ...  
[www.facebook.com/johnmccain](http://www.facebook.com/johnmccain) · [Cached page](#)

# Summaries

- The title is often automatically extracted from document metadata. What about the summaries?
  - This description is crucial.
  - User can identify good/relevant hits based on description.
- Two basic kinds:
  - Static
  - Dynamic
- A **static summary** of a document is always the same, regardless of the query that hit the doc
- A **dynamic summary** is a *query-dependent* attempt to explain why the document was retrieved for the query at hand






# Static summaries

- In typical systems, the static summary is a subset of the document
- Simplest heuristic: the first 50 (or so – this can be varied) words of the document
  - Summary cached at indexing time
- More sophisticated: extract from each document a set of “key” sentences
  - Simple NLP heuristics to score each sentence
  - Summary is made up of top-scoring sentences.
- Most sophisticated: NLP used to synthesize a summary
  - Seldom used in IR; cf. text summarization work

# Dynamic summaries

- Present one or more “windows” within the document that contain several of the query terms
  - “KWIC” snippets: Keyword in Context presentation

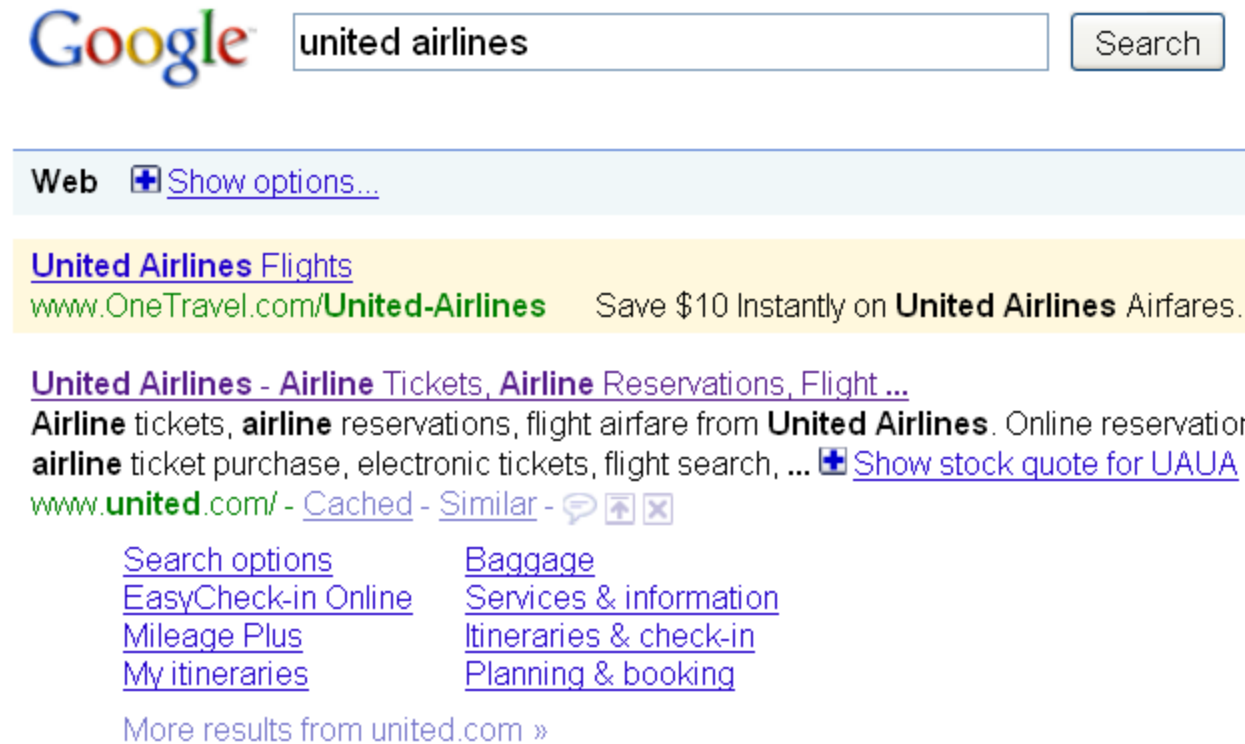
	<input type="text" value="christppher manning"/>	<p><u><a href="#">Christopher Manning, Stanford NLP</a></u>  <b>Christopher Manning</b>, Associate Professor of Computer Science and Linguistics, Stanford University.  <a href="http://nlp.stanford.edu/~manning/">nlp.stanford.edu/~manning/</a> - 12k - <a href="#">Cached</a> - <a href="#">Similar pages</a></p>
	<input type="text" value="christopher manning machine translation"/>	<p><u><a href="#">Christopher Manning, Stanford NLP</a></u>  <b>Christopher Manning</b>, Associate Professor of Computer Science and Linguistics, ... computational semantics, <b>machine translation</b>, grammar induction, ...  <a href="http://nlp.stanford.edu/~manning/">nlp.stanford.edu/~manning/</a> - 12k - <a href="#">Cached</a> - <a href="#">Similar pages</a></p>
	<div> <a href="#">web</a> <a href="#">images</a> <a href="#">video</a> </div> <input type="text" value="christopher manning"/>	<p><u><a href="#">Christopher Manning, Stanford NLP</a></u>  <b>Christopher Manning</b>, Associate Professor of Computer Science and Linguistics, Stanford University ... <b>Chris Manning</b> works on systems and formalisms that can ...  <a href="http://nlp.stanford.edu/~manning/">nlp.stanford.edu/~manning</a> - <a href="#">Cached</a></p>

# Techniques for dynamic summaries

- Find small windows in doc that contain query terms
  - Requires fast window lookup in a document cache
- Score each window wrt query
  - Use various features such as window width, position in document, etc.
  - Combine features through a scoring function – methodology to be covered Nov 12<sup>th</sup>
- Challenges in evaluation: judging summaries
  - Easier to do pairwise comparisons rather than binary relevance assessments

# Quicklinks

- For a *navigational query* such as **united airlines** user's need likely satisfied on [www.united.com](http://www.united.com)
- Quicklinks provide navigational cues on that home page



The screenshot shows a Google search interface. The search bar contains the text "united airlines" and a "Search" button. Below the search bar, there is a "Web" tab and a link to "Show options...". The search results are displayed on a yellow background. The first result is "United Airlines Flights" with a link to "www.OneTravel.com/United-Airlines" and a snippet "Save \$10 Instantly on United Airlines Airfares.". The second result is "United Airlines - Airline Tickets, Airline Reservations, Flight ..." with a snippet "Airline tickets, airline reservations, flight airfare from United Airlines. Online reservation airline ticket purchase, electronic tickets, flight search, ..." and a link to "Show stock quote for UUAU". Below the search results, there are several quicklinks: "Search options", "EasyCheck-in Online", "Baggage", "Services & information", "Mileage Plus", "Itineraries & check-in", "My itineraries", and "Planning & booking". At the bottom, there is a link to "More results from united.com »".

Google united airlines Search

Web [Show options...](#)

**United Airlines Flights**  
[www.OneTravel.com/United-Airlines](http://www.OneTravel.com/United-Airlines) Save \$10 Instantly on **United Airlines** Airfares.

**United Airlines - Airline Tickets, Airline Reservations, Flight ...**  
Airline tickets, airline reservations, flight airfare from **United Airlines**. Online reservation airline ticket purchase, electronic tickets, flight search, ... [Show stock quote for UUAU](#)  
[www.united.com/](http://www.united.com/) - [Cached](#) - [Similar](#) - [Speech](#) [Print](#) [Close](#)

[Search options](#) [Baggage](#)  
[EasyCheck-in Online](#) [Services & information](#)  
[Mileage Plus](#) [Itineraries & check-in](#)  
[My itineraries](#) [Planning & booking](#)

[More results from united.com »](#)

united airlines



Search Pad



SearchScan - On

102,000,000 results for  
united airlines:



Show All



United Air Lines



Wikipedia

Also try: [united airlines reservations](#), [united airlines flight](#), [More...](#)

**United Airlines - Airline Tickets, Airline Reservations ...** (Nasdaq: [UAUA](#))

Official site for **United Airlines**, commercial air carrier transporting people, property, and mail across the U.S. and worldwide.

[www.united.com](#) - 65k - [Cached](#)

[Planning & Booking](#)

[Shop for Flights](#)

[Itineraries & Check-in](#)

[Special Deals](#)

[Mileage Plus](#)

[Flight Status](#)

[Services & Information](#)

[Customer Service](#)

[more results from united.com »](#)

bing

united airlines



UNITED AIRLINES

[United Airline Fleet](#)

[United Airline Schedule](#)

[United Airlines Reservations](#)

[United Airline Jobs](#)

[Reference](#)

ALL RESULTS

[Cheap Flight Tickets](#) · [www.CheapOair.com](#)

CheapOair - The Only Way to Go!! Find Over 18 Million Exclusive Fares.

[Fly United Airlines](#) · [www.OneTravel.com/United-Airline](#)

Save \$10 Instantly on **United Airlines** Flights. Book Now, Hurry!

Best match

[United Airlines - Airline Tickets, Airline Reservations, Flight ...](#)

[www.united.com](#) · Official site

**Airline** tickets, **airline** reservations, flight airfare from **United Airlines**. Online reservations, **airline** ticket purchase, electronic tickets, flight search, fares and availability ...

[Flights](#)

[Redeem miles](#)

[Check In Online](#)

[Children, pets, & assistance](#)

[My itineraries](#)

[Change your travel plans](#)

[Baggage](#)

[Special deals](#)

Customer service 800-864-8331

RELATED SEARCHES

United Airlines [Flight Status](#)

[US Airways](#)

[Continental Airlines](#)

# Alternative results presentations?

- An active area of HCI research
- An alternative: <http://www.searchme.com/> / copies the idea of Apple's Cover Flow for search results
  - (searchme recently went out of business)



# Web Search: Techniques, algorithms and Applications

## Basic Techniques for Web Search

German Rigau <[german.rigau@ehu.es](mailto:german.rigau@ehu.es)>

[Based on slides by Eneko Agirre ...  
and Christopher Manning and Prabhakar Raghavan]

