# **Chapter IR:II**

- II. Architecture of a Search Engine
  - □ Indexing Process
  - Search Process

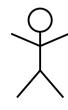
#### Remarks:

□ Software architecture refers to the high level structures of a software system. These structures are needed to reason about the software system. Each structure comprises software elements, relations among them, and properties of both elements and relations.

[Wikipedia]

- □ Software architecture can be specified at various levels of abstraction, also called views. We adopt a high-level functional view, showing what a search engine does, not how it is implemented.
- The software architecture of a search engine must meet two requirements: effectiveness and efficiency. Effectiveness refers to retrieval quality, efficiency to retrieval speed. Other requirements boil down to these two categories. Examples: Scalability demands efficiency; result freshness improves effectiveness and demands efficiency.
- Search engines basically implement two processes, indexing and search, on top of a storage layer. Indexing is a background process to prepare the to-be-searched data for efficient search, and updating it. Search offers a user interface for query submission, and implements query processing, ranking, and result presentation. The storage layer implements a data model for storing documents, index, and logs so that distributed and parallel search are possible.
- □ Compare with Google's early architecture described in "The Anatomy of a Large-scale Hypertextual Web Search Engine" by [Brin and Page 1998]

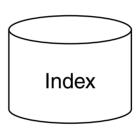






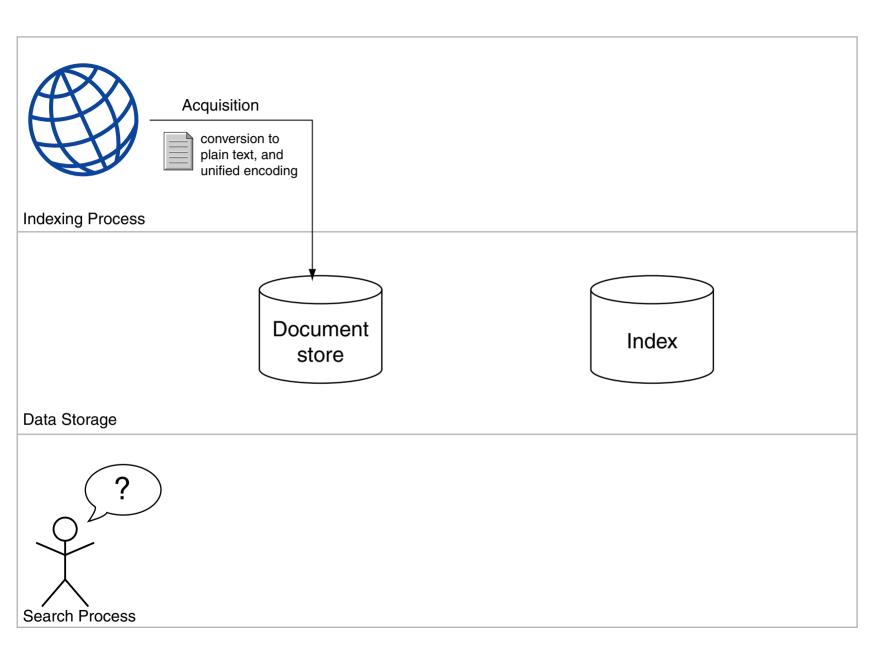






### Data Storage





## Acquisition

In the acquisition step, documents are collected, prepared, and stored.

- Sometimes existing document collections can be used
- More often document collections need to be built

### Components:

- Crawler
- Converter
- Document Store

Acquisition: Crawler

A crawler discovers and acquires documents. Several variants are distinguished:

#### Web crawler

- Exploit hyperlinks to discover web pages
- What are challenges for web crawling?

#### Site crawler

Web crawler restricted to a site (i.e., domain)

## Focused crawler / topical crawler

- Acquires documents matching pre-defined topics, genres, or other criteria; discards others while still exploiting their hyperlinks for discovery
- Requires a document classifier to identify matching documents
- Examples: academic search, news search, business search, job search, etc.

#### Document crawler

- Scans local directories, emails, databases, etc.
- Examples: enterprise search, desktop search

Acquisition: Crawler

## A crawler discovers and acquires documents. Several variants are distinguished:

#### Web crawler

- Exploit hyperlinks to discover web pages
- Exploration policy (e.g., avoid spider traps), duplicate identification (e.g., URL normalization), revisit/update policy, politeness, parallelization

#### Site crawler

Web crawler restricted to a site (i.e., domain)

## Focused crawler / topical crawler

- Acquires documents matching pre-defined topics, genres, or other criteria; discards others while still exploiting their hyperlinks for discovery
- Requires a document classifier to identify matching documents
- Examples: academic search, news search, business search, job search, etc.

#### Document crawler

- Scans local directories, emails, databases, etc.
- Examples: enterprise search, desktop search

#### Remarks:

□ Some web sites inform about updates via web feeds (e.g., using RSS or Atom). For such sites, crawling may boil down to subscribing to the feed, i.e., revisiting the feed to learn about updates.

Acquisition: Converter

In a search engine, the converter unifies documents as follows:

## Reformatting / text extraction

- Documents come in a variety of formats (e.g, HTML, PDF, DOC)
- Subsequent processing steps require unified input format (e.g., plain text)
- Extracting plain text from binary documents is lossy (e.g., formatting is lost)
- Extracting text formatting is relevant, too, for subsequent steps
- The definition of a custom, search engine-specific document format is useful

## Encoding normalization

- Plain text documents come in a variety of encodings (e.g., ASCII, Unicode)
- Subsequent processing steps require unified input encoding (e.g., Unicode)
- Encoding specifications are untrustworthy, encodings must be detected
- Documents encodings are often invalid, they must be repaired
- Errors propagate; when visible in search results, search engine gets blamed

Acquisition: Document Store

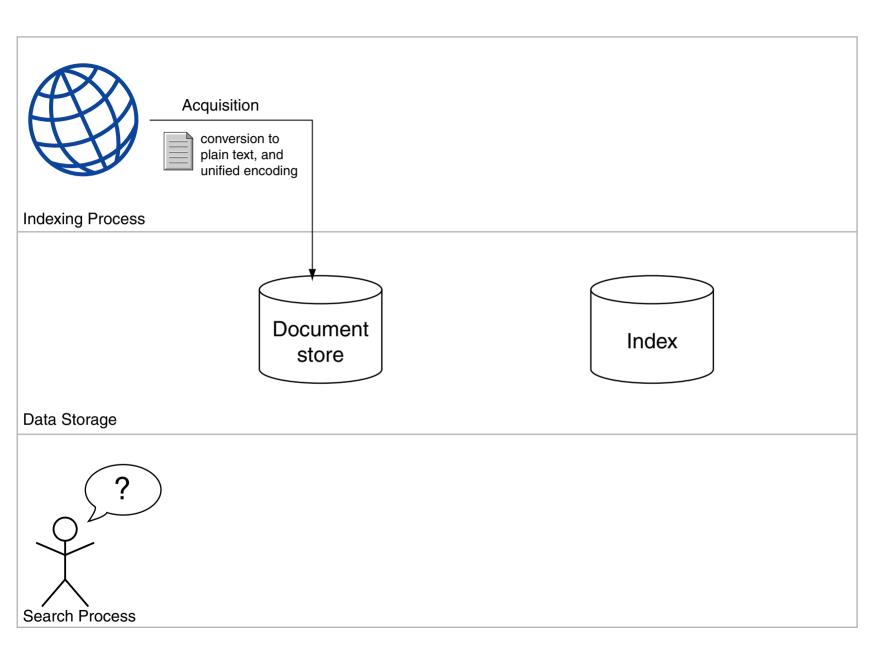
The document store manages the documents acquired:

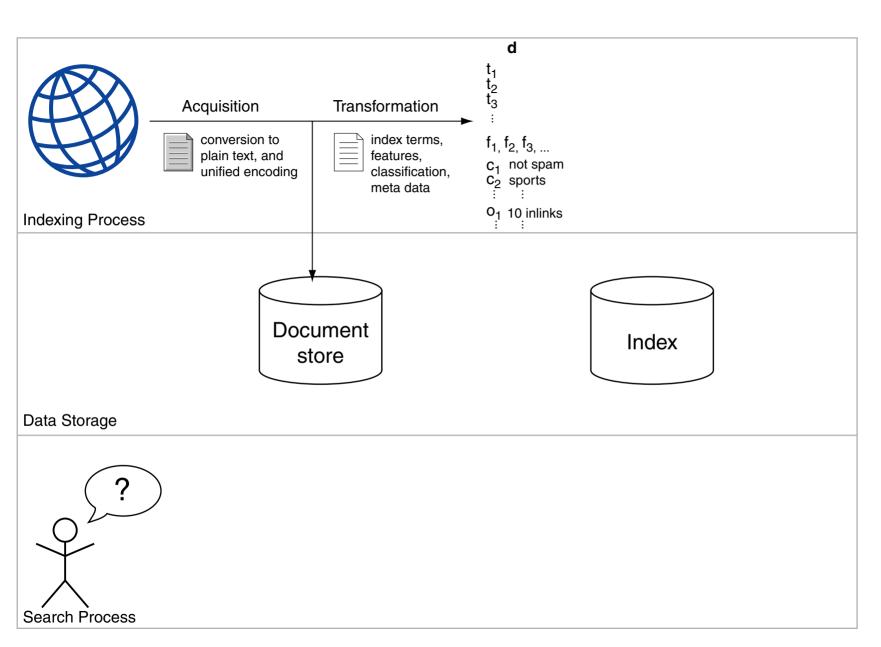
- Original documents are kept to allow for reprocessing as development progresses
  - Why redundant local mirroring of the documents?
- Converted documents are stored for caching (i.e., efficiency)
- Meta data about the documents are stored (e.g., origin, crawl date, etc.)
- Version history of each document may be kept as they get recrawled
- Scale often demands for a distributed document database

Acquisition: Document Store

The document store manages the documents acquired:

- Original documents are kept to allow for reprocessing as development progresses
  - Faster access for reprocessing
  - Original may not always be available
- Converted documents are stored for caching (i.e., efficiency)
- □ Meta data about the documents are stored (e.g., origin, crawl date, etc.)
- Version history of each document may be kept as they get recrawled
- Scale often demands for a distributed document database





#### **Text Transformation**

The transformation step extracts from a document's text the keys by which it can be looked up in the index. Two kinds of keys are distinguished:

- Index terms (terms, for short)
  - Words or phrases of a document's text
  - Their purpose is to represent what a document is about
  - All index terms of all documents combined form the vocabulary

#### Features

- A feature is an individual measurable property of a document
- Their purpose is to represent a document for classification
- Different feature sets are suitable for different classification goals
- Example classification goals: spam, language, genre, . . .
- Both features and derived class labels may be stored

### **Text Transformation**

The transformation step extracts from a document's text the keys by which it can be looked up in the index. Two kinds of keys are distinguished:

- Index terms (terms, for short)
  - Words or phrases of a document's text
  - Their purpose is to represent what a document is about
  - All index terms of all documents combined form the vocabulary

#### Features

- A feature is an individual measurable property of a document
- Their purpose is to represent a document for classification
- Different feature sets are suitable for different classification goals
- Example classification goals: spam, language, genre, . . .
- Both features and derived class labels may be stored

### Components:

□ Segmenter □ Link Extraction

□ Stopping □ Information Extraction

¬ Stemmer / Lemmatizer □ Classification

Text Transformation: Segmenter

Segmentation means breaking down a document into its constituent parts. Two levels of segmentation can be distinguished:

## Page segmentation

- Analysis of the HTML source of a web page with regard to its structure
- Extraction of main content vs. ads, navigation, header, footer, etc.
- Extraction of text structure and text formatting
- Often the HTML of a web page is invalid; it must be processed regardless

### Text segmentation

- Analysis of a web page's plain text with regard to linguistic and structural text units, such as words, sentences, paragraphs
- Tokenization turns a text string into a sequence of token, where a token may be a word or punctuation. Example:
  - Whitespace tokenization: token are separted by white space characters
  - Simple definition of a word: token that is an alphanumeric string
  - Why are these definitions insufficient?
- Lower-cased words are used as index term candidates

Text Transformation: Segmenter

Segmentation means breaking down a document into its constituent parts. Two levels of segmentation can be distinguished:

## Page segmentation

- Analysis of the HTML source of a web page with regard to its structure
- Extraction of main content vs. ads, navigation, header, footer, etc.
- Extraction of text structure and text formatting
- Often the HTML of a web page is invalid; it must be processed regardless

### Text segmentation

- Analysis of a web page's plain text with regard to linguistic and structural text units, such as words, sentences, paragraphs
- Tokenization turns a text string into a sequence of token, where a token may be a word or punctuation. Example:
  - Whitespace tokenization: token are separted by white space characters
  - · Simple definition of a word: token that is an alphanumeric string
  - Punctuation and adjectivization using hyphens not separated by white space
  - Contractions and special characters are neglected
- Lower-cased words are used as index term candidates

Text Transformation: Stopping

Stopping, also stop word removal, discards a selection of words from the set of index terms of a document, or the entire vocabulary. Candidates for stop words:

#### Function words

Words that have little lexical meaning, are ambiguous, serve only grammatical purposes, or specify attitude or mood. Examples: the, of, to, for.

### Frequent words

The most frequently appearing words of a language, or within a collection of documents. Examples: "Wikipedia" appears on every Wikipedia page.

### Domain-specific words

Words that do not discriminate in a given search domain. Examples: "learning" may be ignored in the education domain, regardless its frequency.

Upsides of stopping are reduced index size, faster query processing speed, and potential noise reduction in retrieval. Downsides are that cornercases may not be sufficiently covered. Example: search for "to be or not to be".

Dependent on the retrieval model, stopping may or may not affect retrieval effectiveness. Stopping on text often more conservative than on queries.

Text Transformation: Stemmer / Lemmatizer

Stemming aims at reducing inflected index terms to a common stem. Example: "statistics" should also match "statistic" and "statistical". Two approaches to stemming can be used:

## □ (Heuristic) Stemming

Rule-based affix removal (i.e., suffixes, prefixes, and infixes). Examples: "worker", "megavolt", "un-bloody-likely". Naive heuristic: truncate word at letter 4.

#### Lemmatization

Mapping of a word to its root form, even if it is spelled differently. Example: "saw" and "see".

The upside of stemming/lemmatization is an increased chance of finding a document when it uses different grammar or derived words different from a query.

What are problems related to aggressive stemming?

Text Transformation: Stemmer / Lemmatizer

Stemming aims at reducing inflected index terms to a common stem. Example: "statistics" should also match "statistic" and "statistical". Two approaches to stemming can be used:

## □ (Heuristic) Stemming

Rule-based affix removal (i.e., suffixes, prefixes, and infixes). Examples: "worker", "megavolt", "un-bloody-likely". Naive heuristic: truncate word at letter 4.

#### Lemmatization

Mapping of a word to its root form, even if it is spelled differently. Example: "saw" and "see".

The upside of stemming/lemmatization is an increased chance of finding a document when it uses different grammar or derived words different from a query. Downsides are conflations of unrelated words. Example: "university", "universe", and "universal" are stemmed to "univers" by common stemmers. Stemmed words may not be real ones. Lemmatization requires expensive resources.

Effectiveness depends on language (e.g., Chinese vs. Arabic). Alternative approach: query expansion.

Text Transformation: Link Extraction

### Extraction of links and anchor texts from a document. This serves two purposes:

## Link analysis

Hyperlinks induce a graph among web pages. Link analysis traverses this graph to identify authoritative web pages. Algorithms for this include PageRank and HITS.

## Text augmentation with anchor texts

The text found on a web page may be insufficient to describe its contents. Examples: product pages or pages showing only images. Hyperlinks often enclose text, and the text before and after hyperlinks may justify the link. These anchor texts are added to the text extracted from the linked page to be indexed, too.

Text Transformation: Information Extraction

Information Extraction aims at identifying more complex index terms by means of natural language processing technology (computationally expensive):

### Noun phrases

Phrases which have a noun as its head word, i.e., a noun and any word that modifies it. Examples: "The yellow house is for sale.", "I want a skate board".

#### Named entities

Words or phrases that designate something (e.g., a place, a person, an organization, etc.).

### Coreference resolution

Coreferences, i.e., anaphora and cataphora, are expressions that refer backward or forward in a text, respectively. Resolving *them* is important for text understanding, yet, one of the most difficult problems of natural language processing

#### Relation detection

Extraction of relations between named entities mentioned in the text. Example: "*Bill* lives in the *USA*".

#### Semi-structured information extraction

Extraction of tables, quotes, references, comments, etc.

Text Transformation: Classification

Classification applies machine learning to identify specific types and kinds of documents as well as categorizing them. It is based on the features extracted during transformation. Common classification goals are:

## Spam detection / malware detection

Determines whether a website / web page tries to subvert the search engine's ranking (spam), or to harm its users (malware).

## Language identification

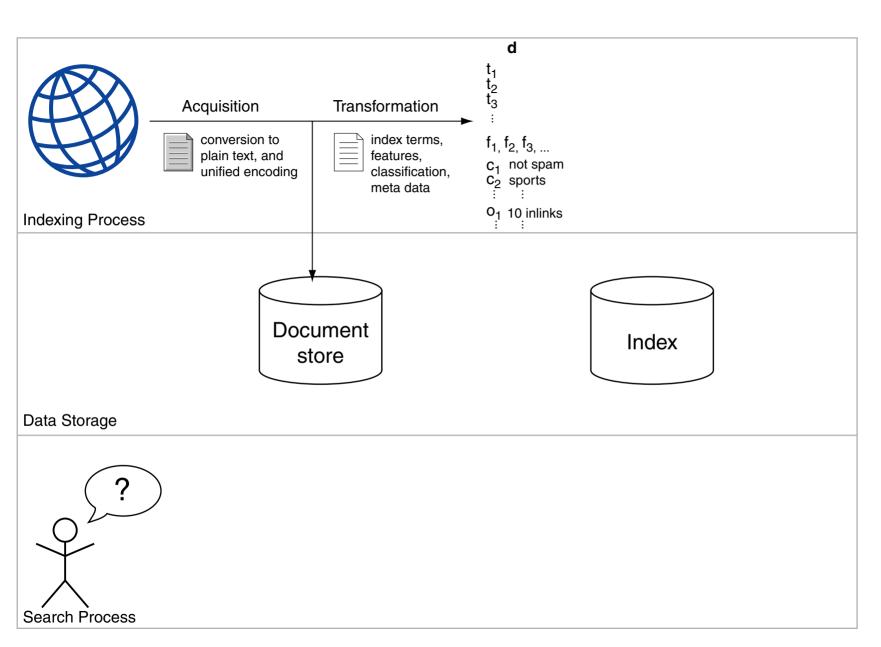
Determines the (main) language of a web page.

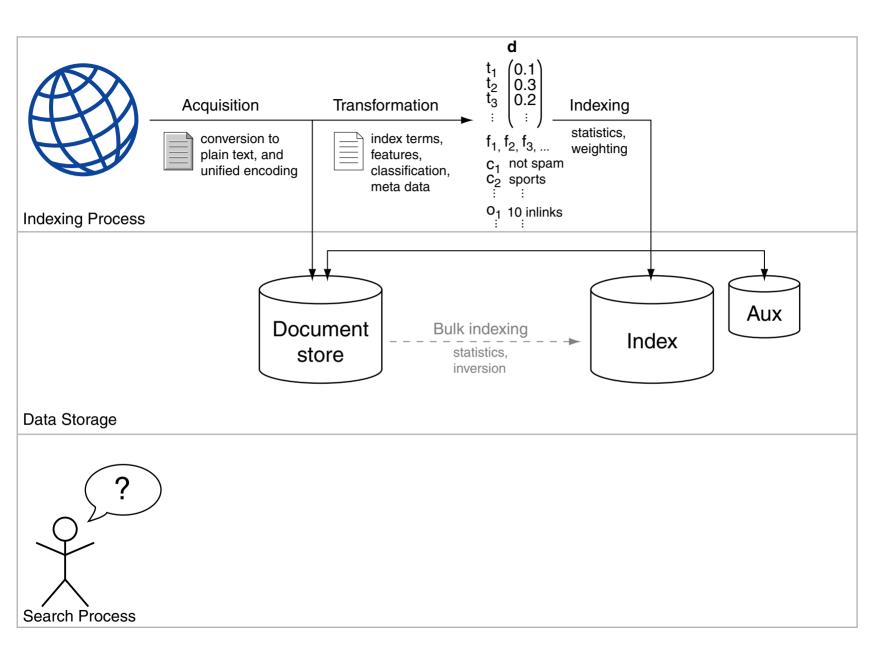
## Topic categorization / cluster analysis

Determines the topic of a document, either by assigning pre-defined subject descriptors (e.g., sports, politics, technology, etc.), covering topics important to users or in a search domain, or by identifying topics using cluster analysis, where identified clusters must be labeled. Topics may overlap and form hierarchies.

## Genre categorization

Determines the genre of a web page (e.g., personal home page, message board, blog, shop, etc.). There is no commonly agreed list of web genres. Genres overlap and form hierarchies, and may depend on search domain.





## Indexing

The indexing step creates the index data structures required for fast retrieval from the transformed documents.

The most commonly used data structure is the inverted index. Other data structures include the suffix array and the signature file.

### Components:

- Document Statistics
- Weighting
- Inversion
- Distribution

Indexing: Document Statistics

Gather and store auxiliary information about documents, required for fast online query processing, e.g., to score and rank documents when including information only available at query time. Such information may include:

- □ Term frequencies per document, topic, and genre
- Document frequencies per term
- Term positions per document
- Document lengths

The data structure used is a key-value store (i.e., a hash map).

Indexing: Weighting

For each index term of a document, calculate a weight indicating its importance with respect to the document's contents, allowing for document scoring with respect to a query. Common term weighting schemes:

- □ Term frequency (tf)
  - Logarithm of the number of occurrences of a term in a document.
- Inverse Document frequency (idf)
  - Document frequency (df): Number of documents comprising a term
  - Logarithm of the number of documents divided by df.
- $\Box$   $tf \cdot idf$

One of the most well-known term weighting schemes in IR.

□ BM25

Similar to  $tf \cdot idf$ , but yields better retrieval performance.

Pre-computing the term weights and storing them in the index or auxiliary data structures speeds up document scoring. Other weighting schemes are calculated based on information about the query.

Indexing: Inversion

Inversion means to change document-term data to term-document data to create an inverted index data structure.

## Bulk indexing

Creates an index by processing all transformed documents. This happens offline; once ready, the currently used index is replaced with the new one.

### Index update

Updates an existing index with new documents as they appear. This happens online, while the index is in use.

Indexing: Distribution

The distribution of an index (also called sharding or partitioning) across machines and data centers facilitates parallel usage.

#### Document distribution

Split collection; smaller indexes for sub-collections on different machines

#### Term distribution

- Split the index for the entire collection by terms
- Different machines serve different terms

### Replication

Copies of (parts of) indexes at multiple sites

□ What are arguments for sharding by documents, by terms, and replication?

Indexing: Distribution

The distribution of an index (also called sharding or partitioning) across machines and data centers facilitates parallel usage.

#### Document distribution

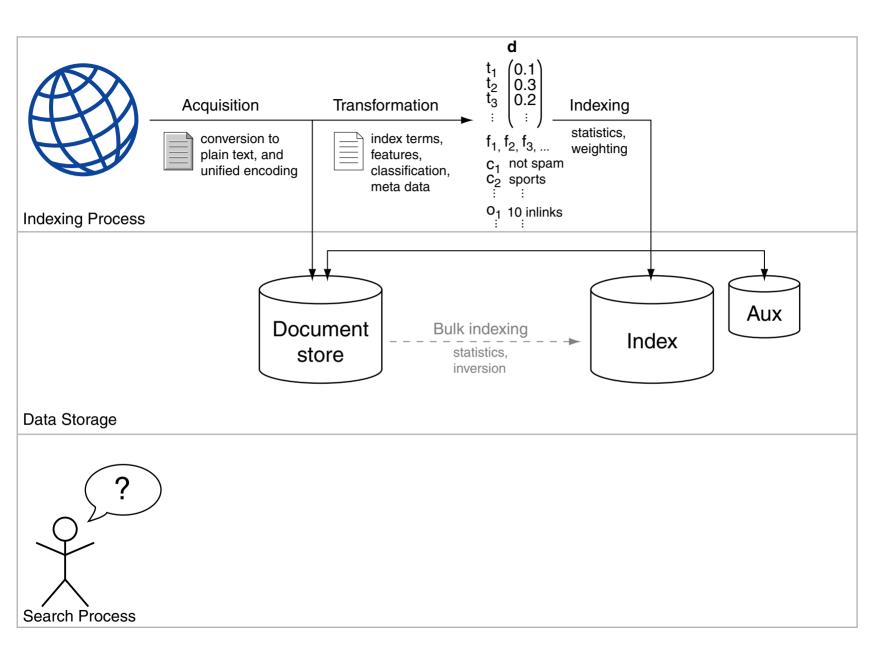
- Split collection; smaller indexes for sub-collections on different machines
- Enables parallelism for indexing and query processing
- Smaller indexes often faster due to caching effects

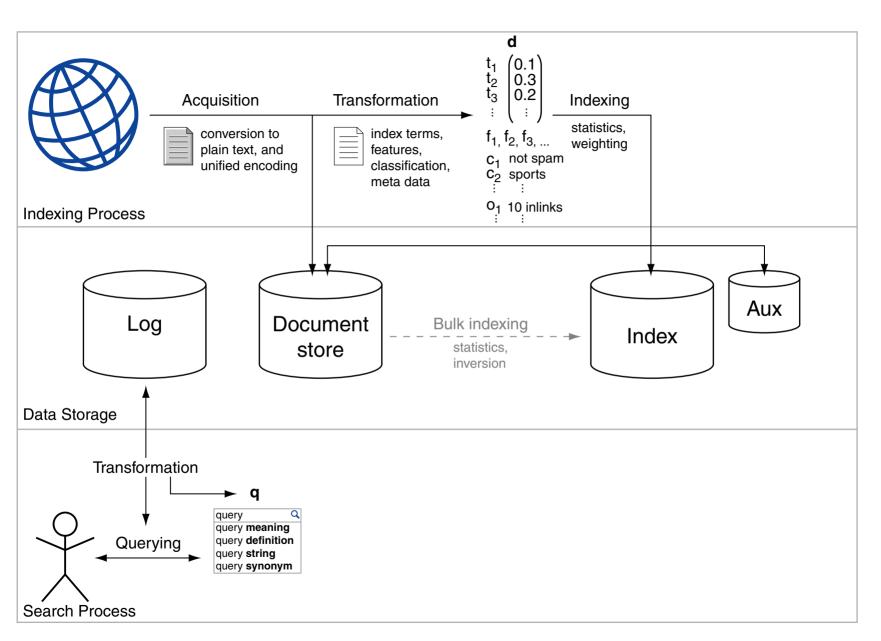
#### Term distribution

- Split the index for the entire collection by terms
- Different machines serve different terms
- Not all machines have to process every query

### Replication

- Copies of (parts of) indexes at multiple sites
- Reduced delays during query processing
- Fault tolerance





### **Search Process**

#### **User Interaction**

User interaction includes the user interfaces offered by a search engine, the usage scenarios envisioned, and their implementation.

#### Basic user interaction:

- Query submission
- Result presentation

#### **User Interaction**

User interaction includes the user interfaces offered by a search engine, the usage scenarios envisioned, and their implementation.

#### Advanced user interaction:

- Query refinement
- Result exploration

#### **User Interaction**

User interaction includes the user interfaces offered by a search engine, the usage scenarios envisioned, and their implementation.

#### Advanced user interaction:

- Query refinement
- Result exploration

## Components:

- Query Language
- Query Transformation
- Results Output

User Interaction: Query Language

The query language defines the syntax and semantics of valid queries. It may include commands to influence the search, so-called query operators.

#### Common query types:

- Structured query
- Keyword query
- Question query
- Query by example

### Common operators:

- □ Boolean operators (AND, OR, NOT (or −))
- What other operators do you know?

User Interaction: Query Language

The query language defines the syntax and semantics of valid queries. It may include commands to influence the search, so-called query operators.

Common query types:	Common operators:
<ul> <li>Structured query</li> </ul>	□ Boolean operators (AND, OR, NOT (or −))
<ul><li>Keyword query</li></ul>	<ul><li>Quotes / phrasal search ("phrase of text")</li></ul>
<ul><li>Question query</li></ul>	□ Field search (title, text, url)
<ul><li>Query by example</li></ul>	□ Wildcards (*, 50100)
	□ Site search (site:example.com)

The most basic form of a query language is the keyword search.

Only about 1% of web queries contain operators [White and Morris 2007], so that web search engines cannot rely on users becoming experts for their language.

Domain-specific search engines often have specialized query languages, allowing for fine-grained control of retrieval behavior.

User Interaction: Query Transformation

The query transformation step maps a query's keywords to index terms, and refines the query in an attempt to better understand the user's intent.

#### Text transformation

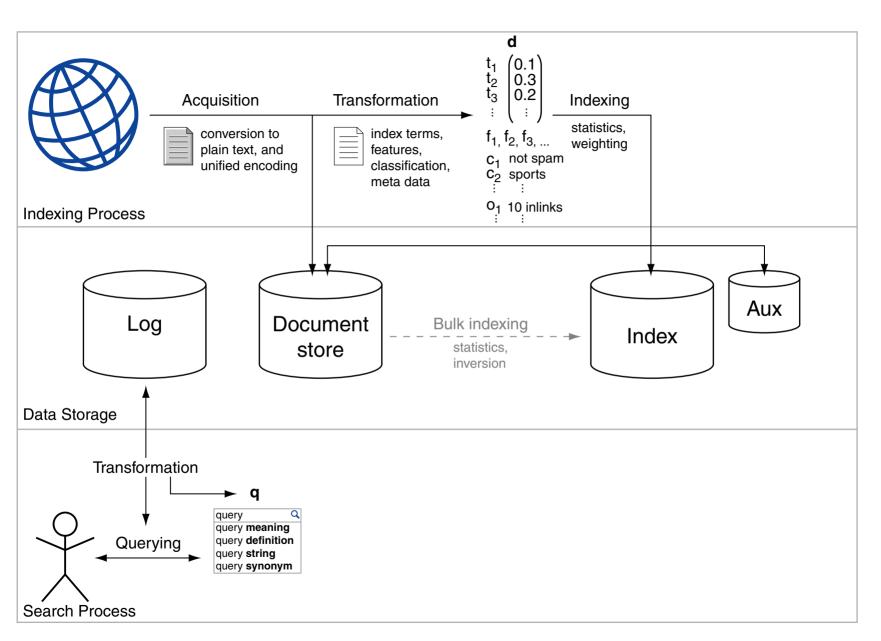
Uses a similar pipeline as the text transformation step for documents, consisting of tokenization, stopping, stemming, etc., to ensure that index terms will match.

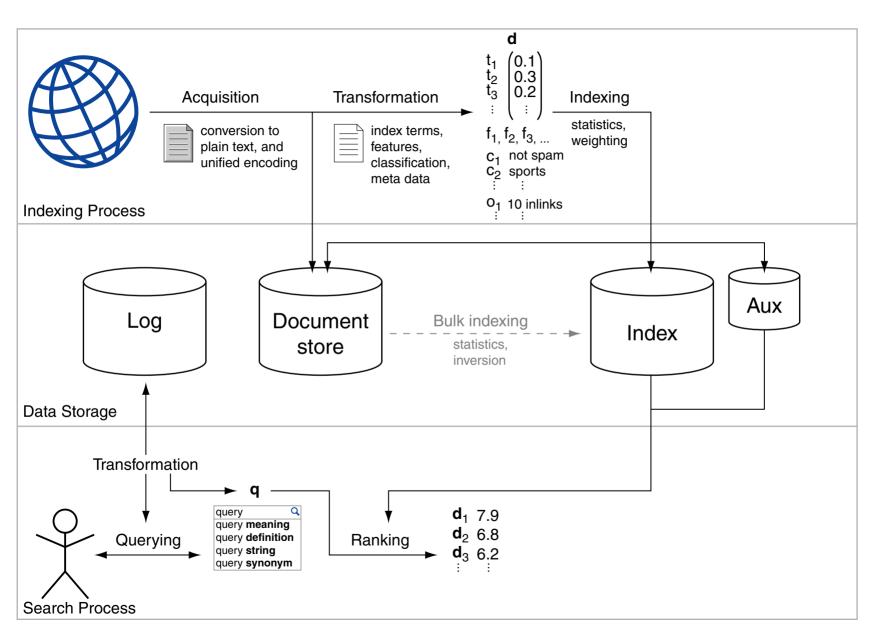
## Spell checking and query suggestion

Gives the user feedback about the query at various degrees of urgency, ranging from small hints ("Did you mean ...") to automatic replacement, dependent on confidence in the suggested alternatives. Query logs are utilized, creating a search log retrieval system.

## Query expansion

Suggestion of word completions, and of additional terms to add to a query, rendering it more specific. Query logs and term co-occurrences in documents are exploited here.





## Ranking

Given a transformed query, the ranking steps scores and ranks the documents indexed with respect to their relevance to the query.

This step marks the keystone of the implementation of the search engine's underlying retrieval model, a theory of how relevance can be quantified.

Retrieval models consists of a function to represent documents and queries, and a function to rank them based on the representations. Document representations are typically pre-computed offline and indexed (e.g., the term weights under  $tf \cdot idf$ ). Rankings are computed online for each query.

### Components:

- Document Scoring
- Efficient Document Scoring
- Distribution

Ranking: Document Scoring

The scoring step quantifies the relevance of the documents indexed to a query.

Ranking: Document Scoring

The scoring step quantifies the relevance of the documents indexed to a query.

Let  $t \in V$  denote a term t from the vocabulary V of index terms, and let  $\omega_X : V \times X \to \mathbf{R}$  denote a term weighting function, where X may be sets of documents D or queries Q. Then the most basic form of scoring a document is

$$\sum_{t \in V} \omega_q(t) \cdot \omega_d(t),$$

where  $\omega_q(t)$  and  $\omega_d(t)$  are term weights indicating the importance of t for the query  $q \in Q$  and the document  $d \in D$ , respectively.

Ranking: Document Scoring

The scoring step quantifies the relevance of the documents indexed to a query.

Let  $t \in V$  denote a term t from the vocabulary V of index terms, and let  $\omega_X : V \times X \to \mathbf{R}$  denote a term weighting function, where X may be sets of documents D or queries Q. Then the most basic form of scoring a document is

$$\sum_{t \in V} \omega_q(t) \cdot \omega_d(t),$$

where  $\omega_q(t)$  and  $\omega_d(t)$  are term weights indicating the importance of t for the query  $q \in Q$  and the document  $d \in D$ , respectively.

#### Observations:

- $\Box$  The term weights  $\omega_d(t)$  have been pre-computed and indexed.
- $\Box$  The term weights  $\omega_q(t)$  must be computed on the fly.
- □ A term t may have importance, and hence non-zero weights, for query or document despite not occurring in them. Example: synonyms.
- $\Box$  The majority of terms from V will have insignificant importance to both.

## Ranking: Efficient Document Scoring

Computing document scores requires index access. Different access strategies govern what can be accomplished and how the data has to be organized. The two most salient strategies are as follows:

## Document-at-a-time scoring

- Precondition: a total order of documents in the index's postings lists is enforced (e.g., ordering criterion document ID or rather document quality).
- Postlists of a query's terms are traversed in parallel to score one document at a time.
- Each document's score is instantly complete, but the ranking only at the end.
- Concurrent disk IO overhead increases with query length.

## Term-at-a-time scoring

- Traverse postlists one a time (e.g., term ordering criterion frequency or importance).
- Maintain temporary query postlist, containing candidate documents.
- As each document's score accumulates, an approximate ranking becomes available.
- More main memory required for maintaining temporary postlist.
- Safe and unsafe optimization exist, e.g., to stop the search early.

Ranking: Distribution

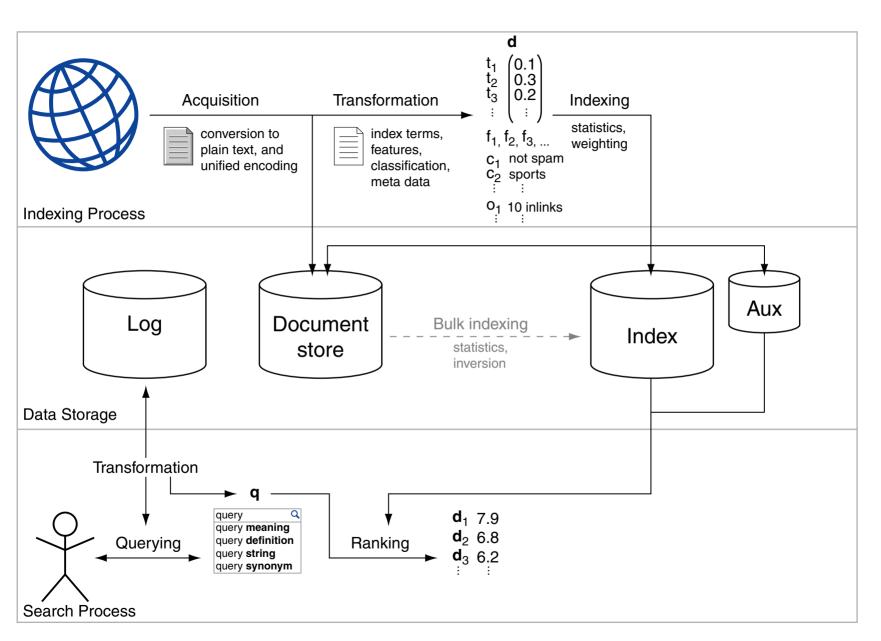
The distribution of query processing depends on that of the index.

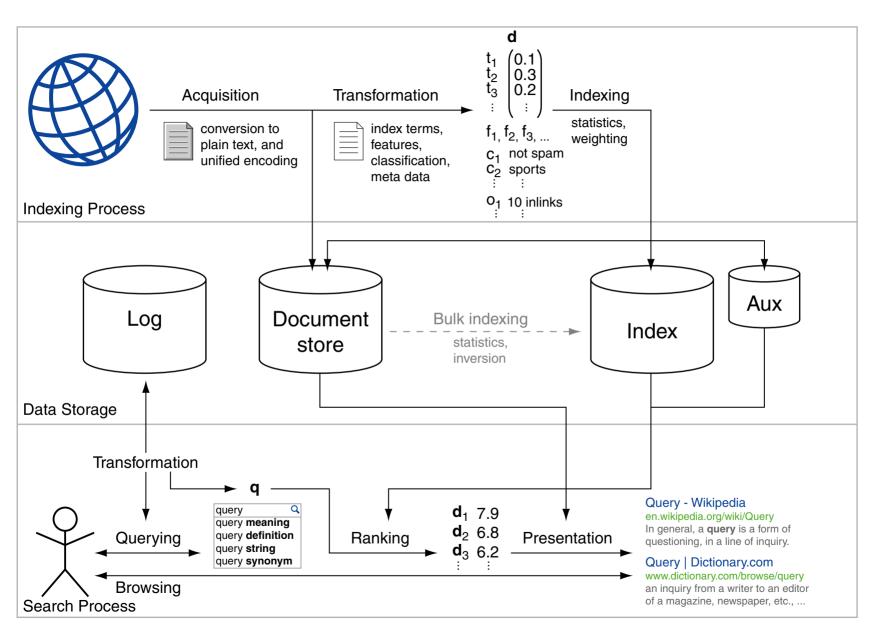
## Query broker / load balancer

Decides which shard and which replicated copies to access. Receives and merges results.

#### Cache

Keeps frequently used data close at hand (e.g., in main memory) to reduce latency. Caches may include indexes containing important documents only held in main memory, precomputed search results, temporary postlists, parts of postlists, and the optimized usage of caching hierarchy from operating system to hardware caches.





User Interaction: Results Output

The results output step compiles the web page displayed to the users. Several additional retrieval steps are required to compose the web page:

## Snippet generation

Accesses the original web page and extracts sentences and phrases that summarize it, dependent on the query.

## Query term highlighting

Preprocesses the snippet to highlight words from the query, regardless their inflection.

## Clustering

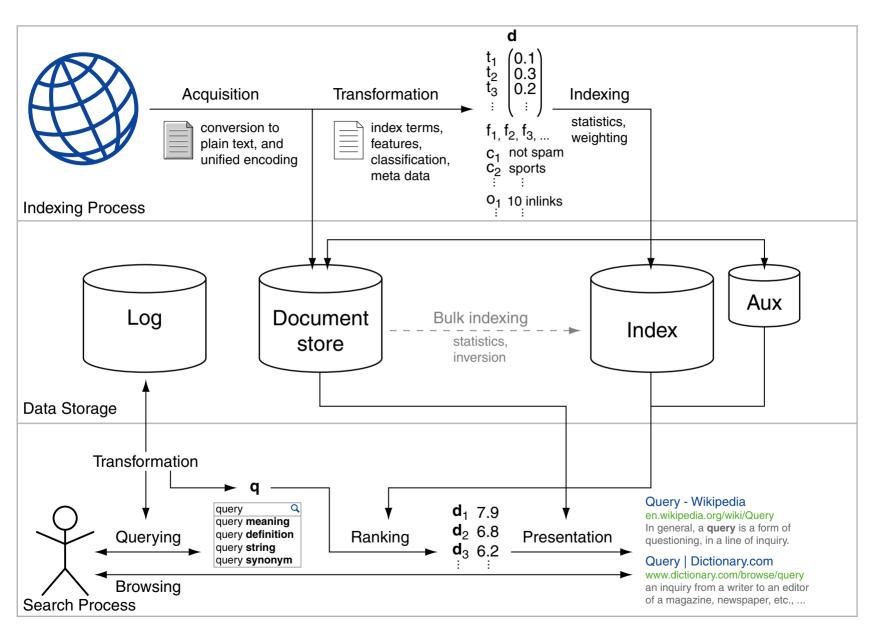
Optionally clusters the set of results to give a more diverse set of results.

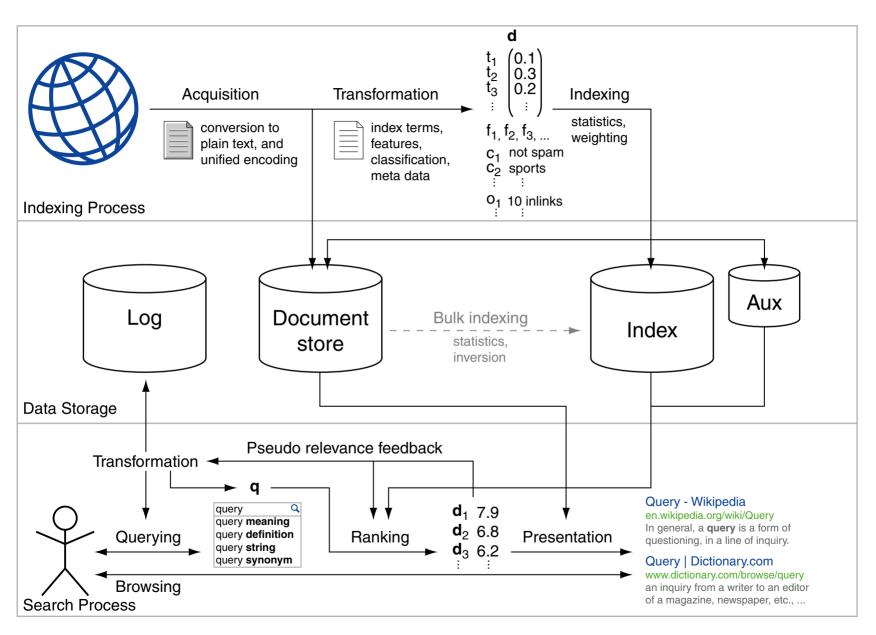
#### Ad retrieval

Accesses another search engine specifically tailored to the retrieval of ads relevant to a query from all ads offered by advertisement partners. This runs in parallel to the retrieval of organic search results.

#### Other oneboxes

Determines whether other specialized search engines can supply relevant results. Ranks the oneboxes into the search results as per their importance compared to organic web results.





User Interaction: Query Transformation (continued)

The query transformation step maps a query's keywords to index terms, and refines the query in an attempt to better understand the user's intent.

#### Text transformation

Uses a similar pipeline as the text transformation step for documents, consisting of tokenization, stopping, stemming, etc., to ensure that index terms will match.

## Spell checking and query suggestion

Gives the user feedback about the query at various degrees of urgency, ranging from small hints ("Did you mean ...") to automatic replacement, dependent on confidence in the suggested alternatives. Query logs are utilized, creating a search log retrieval system.

## Query expansion and relevance feedback

Suggestion of word completions, and of additional terms to add to a query, rendering it more specific. Query logs and term co-occurrences in documents are exploited here.

Pseudo relevance feedback adds terms to a query extracted from the top search results. What problem do you see with pseudo relevance feedback?

User Interaction: Query Transformation (continued)

The query transformation step maps a query's keywords to index terms, and refines the query in an attempt to better understand the user's intent.

#### Text transformation

Uses a similar pipeline as the text transformation step for documents, consisting of tokenization, stopping, stemming, etc., to ensure that index terms will match.

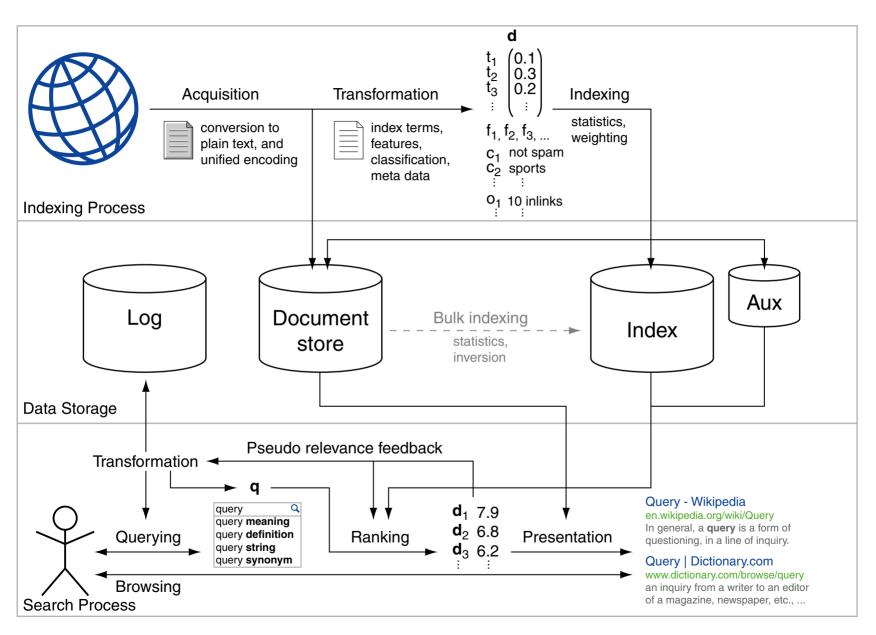
## Spell checking and query suggestion

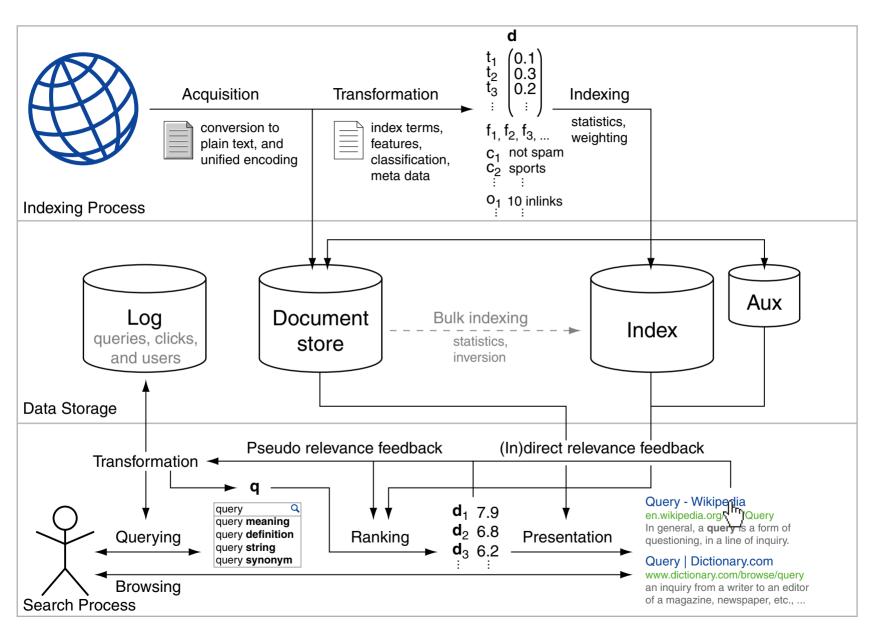
Gives the user feedback about the query at various degrees of urgency, ranging from small hints ("Did you mean ...") to automatic replacement, dependent on confidence in the suggested alternatives. Query logs are utilized, creating a search log retrieval system.

## Query expansion and relevance feedback

Suggestion of word completions, and of additional terms to add to a query, rendering it more specific. Query logs and term co-occurrences in documents are exploited here.

Pseudo relevance feedback adds terms to a query extracted from the top search results. The topic of the top search results is reinforced, which may differ from the user's intent.





## Logging

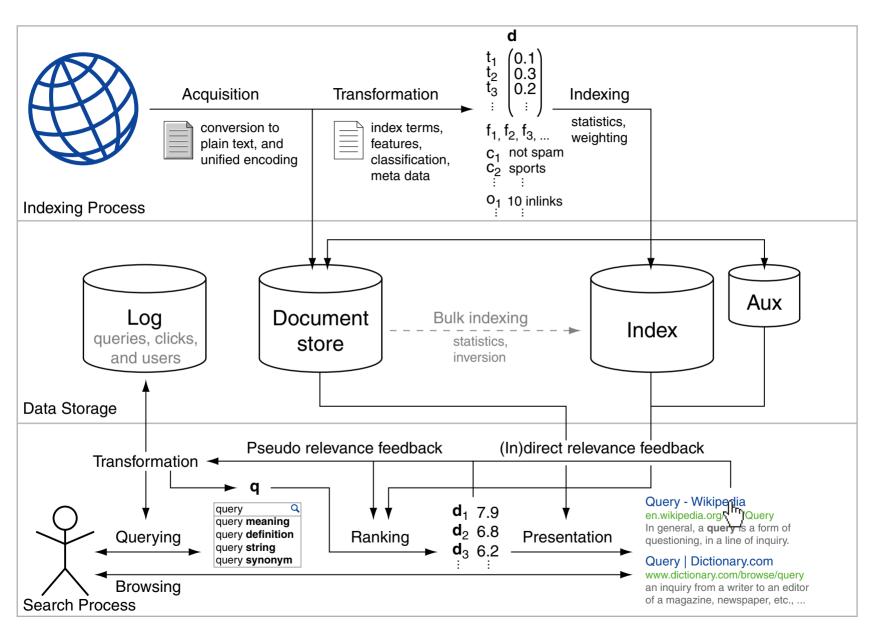
Log all users' activities, specifically queries and interactions with search results.

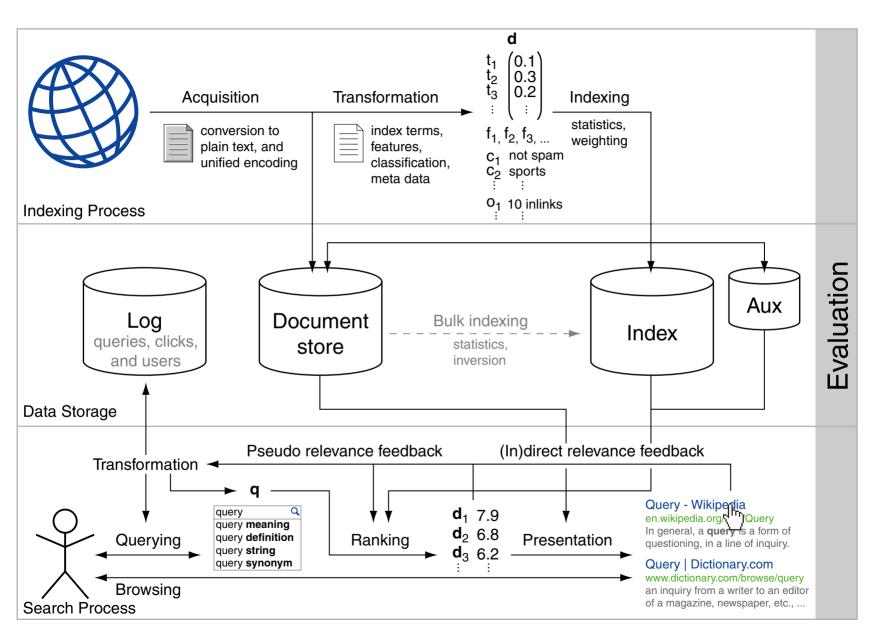
- Query understanding and refinement
  - The log serves the query transformation step.
- □ (In)direct relevance feedback

The log informs the query transformation step and document ranking: direct feedback (applied less often) asks the user which documents returned are relevant, using this information to refine the query and to rerank all documents. Indirect feedback exploits other user's behavior on similar queries to refine the query and train the ranking algorithm to rank better in the future. For the latter, result clicks and estimated dwell times are analyzed.

- Personalization and ad targeting
  - User profiles allow for tailoring search results and ads to the user's context and interests.
- User experience analysis and optimization
  - The user behavior on the search engine's web pages allows for conclusions about its efficacy in supporting the user. Example: colorization of search results has been optimized.

Logs are the most valuable data a search engine collects. A search engine specifically for search in search logs is typically required.





### **Evaluation**

#### Overview

The evaluation of a search engine addresses the analysis of its effectiveness and efficiency.

## Ranking analysis

- Definition of goals
- Acquisition of relevance judgments for query-document pairs (e.g., via crowdsourcing)
- Measurement theory (e.g., strong emphasis on top results common in web search)
- Log analysis of recorded search behavior
- User studies and A/B tests

## User experience analysis

- Definition of goals (e.g., usability, user satisfaction, etc.)
- Log analysis of recorded user behavior
- User studies and A/B testing

## Performance analysis

- Definition of goals (e.g., throughput, response time, etc.)
- Log analysis of recorded system behavior
- Lab experiments and simulation

#### Remarks:

Evaluation is a systematic determination of a subject's merit, worth, and significance, using criteria governed by a set of standards. It can assist to ascertain the degree of achievement or value in regard to the aim and objectives sought after. The primary purpose of evaluation, in addition to gaining insight into prior or existing initiatives, is to enable reflection and assist in the identification of future change.

[Wikipedia]

**Architecture of a Search Engine** 

