# Information Retrieval

**Overview**

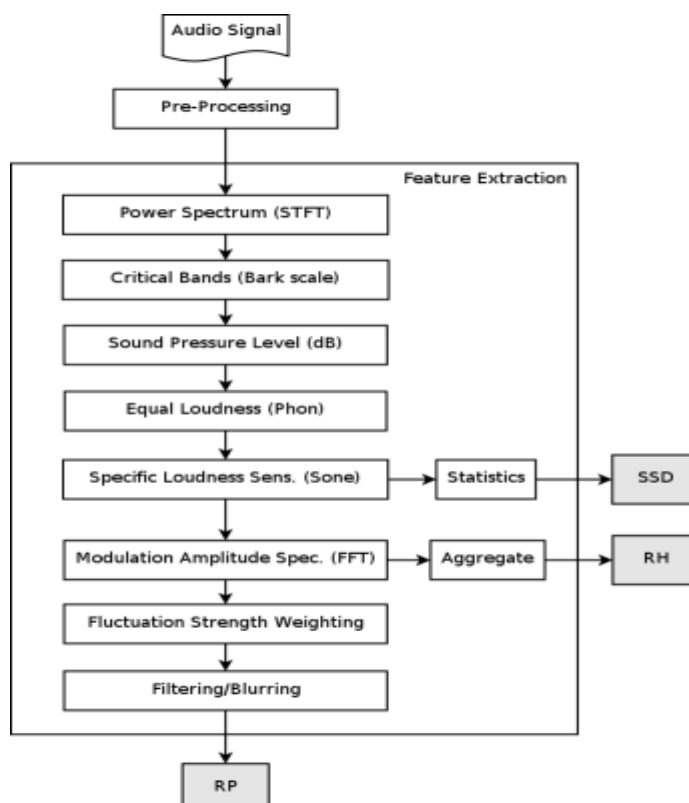|  | **IE (Data base)** | **IR (set of result docs)** |
|---|---|---|
| **data/documents** | structured | unstructured |
| **attribute semantics** | defined | ambiguous |
| **queries** | well defined | Can be free text/arbitrary |
| **retrieval** | exact | imprecise |

**Procedure IR:**

- Pre-Processing
  - Collection Cleansing
    - stop-words
    - formatting information
    - empty/ill-formated docs
  - identify relevant doc parts for indexing (text, image, meta data, ...)
  - stemming (e.g. Porter Stemmer)
  - stop-word removal (manually created stop-word list)
- Indexing (Bag of words)
  - select type of terms used (feature set selection)
    - n-grams
    - words (word stems)
    - word co-occurences (word n-grams) $\rightarrow$ detect phrases
    - concepts (Date, Person, Company, ...) $\rightarrow$ NLP
  - and weight the terms $\rightarrow$ term weighting (feature extraction)
    - df (document frequence) of term to select terms
    - tf (term frequence) inside a document
    - tfidf (tf inverse df) = tf/df oder tfidf = tf/(ln(N/df))
    - 
    - Zipf's law (relates term frequency to rank)
    - Heap's law (Predict number of distinct terms)
    - 
- Retrieval

- find documents satisfying query
- Retrieval Models
    - Boolean model (Exact-Match)
        - canonicalization until result is satisfactory
        - detailed knowledge of doc domain needed
    - Vector Space Model (Best-Match)
        - normalized high-dimensional feature space (indexing)
        - query = vector
        - result is docs that are closed to high-dimensional query vector
            - similarity by L1, L2, minkovsky, or cosine similarity (distance measure)
    - Probabilistic Model
        - 2-class classification (relevant or not)
        - e.g. Bayes statistics
- Relevance Feedback
    - iterative interactive retrieval to refine query (e.g. add terms)
        - manual refinement (blind – user define relevant result docs)
        - semi-automatic refinement
        - automatic refinement (see ATC) → e.g. rocchio feedback: add terms from relevant and substract terms form irrelevant docs
- ATC – Automatic Text Classification
    - Machine Learning
        - knn-classifiers
        - decision trees
        - rocchio
        - naive bayes
        - support vector machines
        - clustering (SOM, etc.)
    - feature space dimensionality reduction
        - feature selection vs. Extraction
        - local vs. Global
- Evaluation
    - Contingency table
    - Measures
        - precesion
        - recall
        - accuracy/error

- F-measure

**Music IR**

- What is Music?
  - Sound
    - Nyquist sampling theorem
    - lossless/lossy sound formats
    - PCM (Pulse Code Modulation) → digital representation of analog signal
    - MIDI (Musical Instrument Digital Interface)
    - Scores/Sheet-Music (by hand → many styles, Printed, MusicXML, (e.g. Lily Pond), ...)
  - Text
  - Community Data
  - Images/Videos
- Web Music Retrieval
  - music search engines
  - centralized/de-centralized/hybrid P2P
- Audio features
  - MPEG7-Standard Features
  - Marsyas System
  - Rhythm Patterns/Rhythm Histograms/Statistical Spectrum Descriptor

- Evaluation/Benchmarking of Music classification (clustering/similarity based) and Music IR
  - E.g. MIREX (besides many others)
- Application Example
  - SOM Player using clustering
  - online (client-based) audio feature extraction
  - audio segmentation - lead-in, verse, chorus, etc. (→ structure e.g.: ABCBDBAC'A) & e.g. k-means clustering
  - chord detection
  - instrument separation using template matching → artificial music

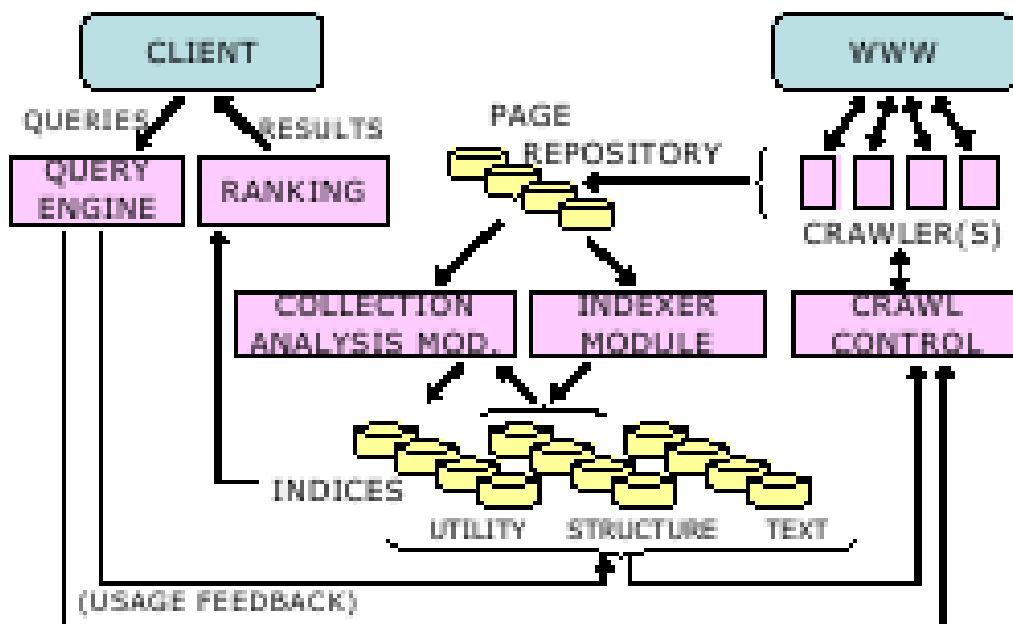**Information Extraction (IE)**

- Application Areas
  - Database population → Filling slots in a database from sub-segments of text
  - Ontology population/evolution
  - Automatic summarization - abstract/extract main ideas in less space
    - Input
      - singe- VS multi-document summarization
      - mono- VS multi- VS cross-lingual
      - text VS multimedia
    - Purpose
      - informative VS indicative(user-interaction)
      - generic VS user-orientated(based on user query)
      - domain-specific VS general (portable to all domains)
    - Output – Quality is crucial
      - Extract(sentences, paragraphs, ...)
        - machine learning/language processing (e.g. tree generation)
        - Edmundsonian Paradigm (ranking of sentences) → can include post-processing to eliminate incoherences (anaphora, semantic gaps, etc..)
      - Abstract (scripts & simple generation VS MUC-like concepts & NLG)
        - prior knowledge (semantic structure known → ontology)
        - clustering
    - Evaluation
      - intrinistic → summary's purpose/task is neglected
        - sentence integrity (anaphor without reference)
        - readability
        - fidelity

- human „gold" summary → precision & recall & etc.
  - extrinsic → account for purpose/task
    - relevance (compare to topic(s) assigned by human judges)
    - comprehensibility (judges are asked questions on studied text)
  - Question answering (natural language answering)
    - Dimensions
      - nature of information (DB ↔ free text)
      - nature of questions (facts ↔ opinions)
      - nature of answers (extracted ↔ [NLP-]generated)
      - nature of technique (linguistically correct/sophisticated ↔ linguistically uninformed (e.g. n-grams)
    - Components
      - Analyze question
      - Gather Information
      - Distill Answer
      - Sanity check
      - Present Answer
  - Web Search (mostly ranked list, no precise answer)
    - Components
      - Crawler (Robots, Spiders)
        - Page selection based on
          - Factors
            - Coverage
            - Quality (indexing good pages)
            - Efficiency (no duplicates)
            - Etiquette (minimize overloaded server loads)
            - Freshness/Age (page life-time)
              - Batch Mode Crawler (periodic update)
              - Steady Crawler (incremental update)
                - keep local collection fresh
                - continuously improve collection's value
              - determined with fixed/variable frequency (based on page's rate of change)
              - in-place update VS shadowing
          - Mathematical Model
            - Importance Metrics

- - - Interest-Driven
    - Popularity-Driven
    - Location-Driven
  - Quality Metrics
    - performance of a crawler is described
  - Ordering Metrics
- Page repository (scalable storage system – local copy of the web)
  - Interface to
    - Crawler
    - Indexer Module
    - Query Engine
  - Storage Manager
    - distribute Pages over available storage nodes
      - uniform distribution policy
      - hash distribution policy
      - log-structured policy (B*-tree index)
      - 
    - handles updates
    - handles different access modes (stream, random access, ...)
- Indexer (uses statistician during merging OR during flushing)
  - Indices Types
    - Utility Index (Search engine specific info → speed up!) - e.g. rank, site index, ...
    - structure/link index (graph-modelled → VERY large graph representing the links)
    - text index=inverted file (identify and select pages) - e.g. rate of change, anchors, headings, ...
      - → is a mapping of content to its location!!!
        - local
        - global
  - Parallel Processing
  - 
- Ranking
  - Extended Boolean Models
    - tf, headings, titles, keywords,...
    - idf, word count,...

- o Content-based (Vector or Probabilistic Model)
- o human annotation
- o Factors
  - ad-hoc factors (porn filter...)
  - popularity
  - text anchor
  - LINK ANALYSIS (find pages with high authority (HUBS) →
    Assumption: good pages link to good pages)
    - e.g. Hypertext Induced Topic Selection(HITS) Algorithm:
      - o generate query-independent sub- graph of THE web graph
      - o recursively calculate hubs and authorities –> refine graph
  - Bibliometric law (often cited articles hae high scientific value)
- o Problems
  - Rank sink
  - Rank leak
- Architecture



- 3<sup>rd</sup> Generation Web Search Engines
  - semantic analysis (what does it mean?)
  - determine user context rather than analyze query
    - o user location
    - o previous queries
    - o user profile
    - o spell checking

- o query suggestion
- General Procedure
  - o segmentation (select relevant words/terms/phrases)
  - o classification (noun, verb, concept, ...)
  - o clustering (group classified data by co-referencing detection – e.g. group = „Gerhardt works Apple")
  - o association (fill groups into DB)
- Architecture
  - o Tokenization (for text sectioning and filtering → see indexing)
  - o Lexicon and Morphology (for e.g. maximum entropy POS Tagging) → Named-Entity Recognition (e.g. „Apple Computer Inc.", „Sepp Maier")
    - ■ Combining Morphemes:
      - • free morphemes → stand alone words
      - • bound morphemes → 'tion' ind 'creation
      - • inflectional morphemes → big, big-'ger', big-'gest'
      - • derivational morphemes → verb + 'ment' == noun
    - ■ choice of Morphology or Lexicon
      - • language dependent
      - • domain-specific (medicine, chemistry, literature, ...)
  - o Parsing
    - ■ find grammatical structure (also phrases, etc.)
    - ■ and for Co-reference resolution:
      - • name-aliases
      - • pronoun-antecedents
      - • using definite description like ontology
  - o Domain-Specific analysis (to merge partial results using detected co-referencing)
    - ■ using templates consisting of slots (ontologies
    - ■ Approaches
      - • atomic approach
        - o intelligent guessing →  high recall – low precision
        - o precision improved by filering
      - • molecular approach (more popular)
        - o small amount of highly reliable rules try to match all arguments to pattern/template/ontology → high precision – low recall
        - o iteratively generalizes rules to cover other patterns/templates → Over-generation possible with leads to lower precision and higher recall
- Wrappers

- o   rigorous (unified format, complete info)
- o   semi-rigorous (unified format, incomplete info)
- o   semi-relaxed  (non-unified format, complete info)
- o   relaxed (non-unified format, incomplete info)
- Knowledge Eningeering
  - o   pre-processing (segmentation, filtering, see indexing)
  - o   Analysis (parsing, semantic interpretation – e.g. through co-references)
  - o   post-processing (chose templates/ontologies and categories to map)
  - o   Example models:
    - FASTUS (Finite State Automation Text Understanding System)
    - GENLTOOLSET
    - PLUM (Probabilistic Language Understanding Model)
    - PROTEUS
- Machine Learning
  - o   Supervised learning
    - Autoslog
      - extract dictionary of concepts
      - single slot extraction
    - PALKA
      - new rules are generalized
      - Existing rules are specialized
    - WHISK
      - REGEXP in top-down induction
    - RAPIER (Robust Automated Production of Extraction Rules)
      - Input: already filled templates
      - Output: Pattern-Matching rules
      - bottom-up
    - GATE (General Architecture for Text Engineering)
    - WIEN (Wrapper Induction Environment)
    - bottom-up
    - Lixto
      - Supervised Wrapper Generation Program
  - o   Semi-supervised learning
    - Bootstrapping
  - o   Unsupervised learning