

Advanced Human Machine Interaction

Personalised information retrieval over the Internet

Alexandre Pauchet

alexandre.pauchet@insa-rouen.fr - BO.B.RC18



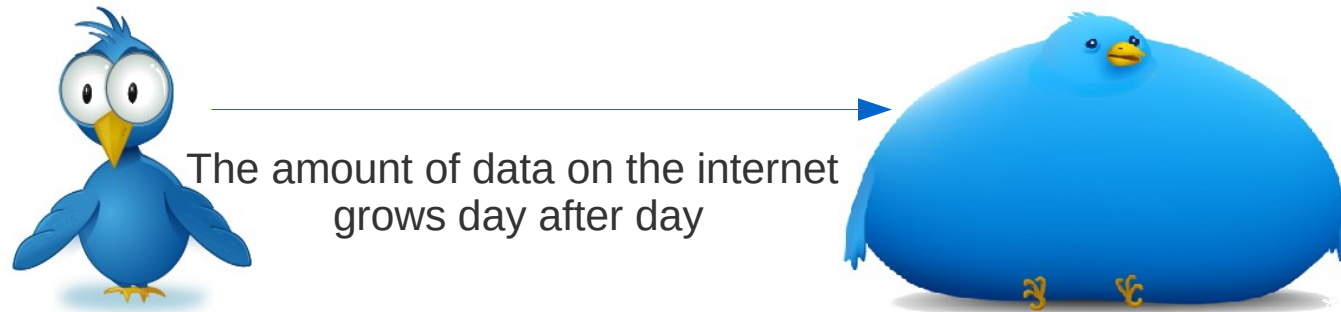
Normandie Université



Definitions

- **Information Retrieval (IR):** finding documents relevant to a user query from a huge database.
- **Document:** text, piece of text, webpage, image, video... In IR, a document corresponds to any item that can answer a user query.
- **Query:** formalisation of a user's need of information. Often, it is a conjunction of keywords.
- **Relevance:** metric evaluating the match between a (found) document and a user query.

Internet as an information source



- Content in **constant growth**
- **Heterogeneous** content (texts, videos, pictures...)
- **Dynamic** content (new pages, dynamic pages : blogs/news/forums...)
- **Few structured content** (metadata, hypermedias, tagging semantics...), insufficiently exploited

Internet = Open content

- **Free and simple posting:** unmonitored content (personal webpages, blogs, wiki, forums)
 - Fake information
 - Verified information then modified
 - Information validated by stakeholders
 - Popularity does not mean truth
- => Information relevance and information validity must be evaluated depending on a need

Information retrieval over the Internet

- User satisfaction difficult to measure
=> Exploitation of relevance metrics
- Page “qualities” being very different, relevance should depend on it
- **Quality of research user interface depends on**
 - Speed
 - Index size
 - Error robustness (approximations, bad wording, ambiguities...)
 - Offered "services" (e.g. Google...)

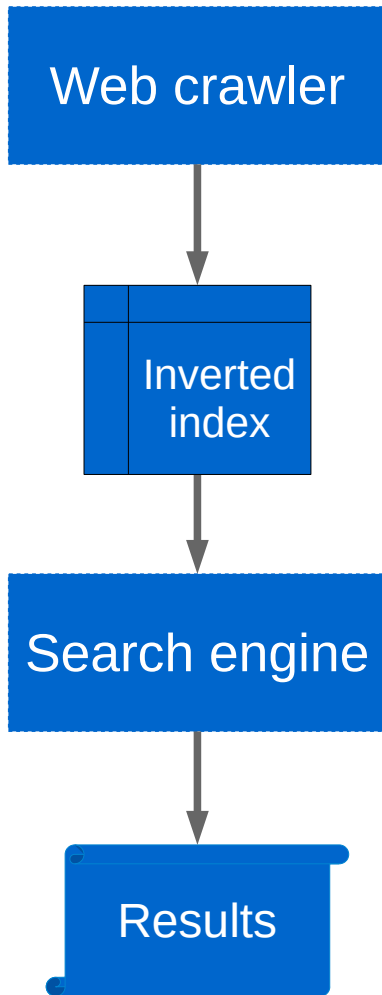
History of IR over the Internet

- **First generation** (~ 1995 : Altavista, Excite, Lycos, etc.) :
 - Indexing from information on the page
- **Second generation** (~ 1998 : Google and others) :
 - Indexing from web structure (entry links, text of links, etc.)
- **Third generation** (~2023) :
 - Answer user's need
 - Semantic and context-specific analysis
 - Help to user: HMI, several languages, autocompletion and spell checking, suggestion of queries...
- **Fourth generation** (ongoing...) :
 - Retrieval-Augmented Generation (RAG)

HMI and IR

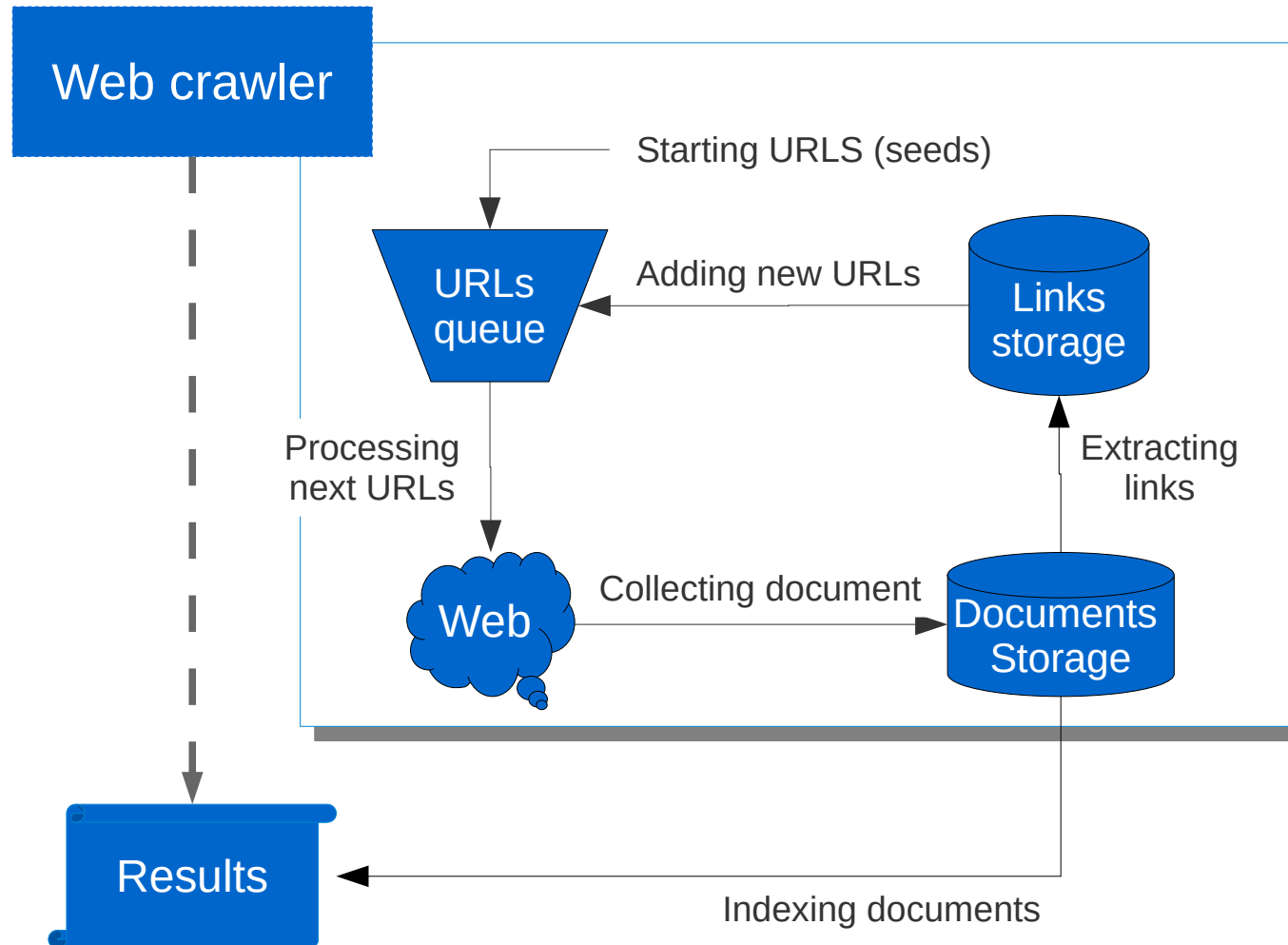
- Representation of documents (index)
- Representation of user's needs
 - Long term needs (profile, topics of interest)
 - Short term needs (context-specific search)
- User interaction ↔ "Poor" IR systems
 - Little or no access to context
 - Limited query (keywords)
 - Document index not personalized
 - Relevance feedback difficult to catch

IR: basic approach



- Crawlers browse websites and index their content
 - Frequency depends on search engines
 - A page can ask to not be indexed (robots.txt), or to not follow links ("nofollow")
- Index is a summary at time T of webpages
- Results of a search engine are ordered by relevance towards a user request

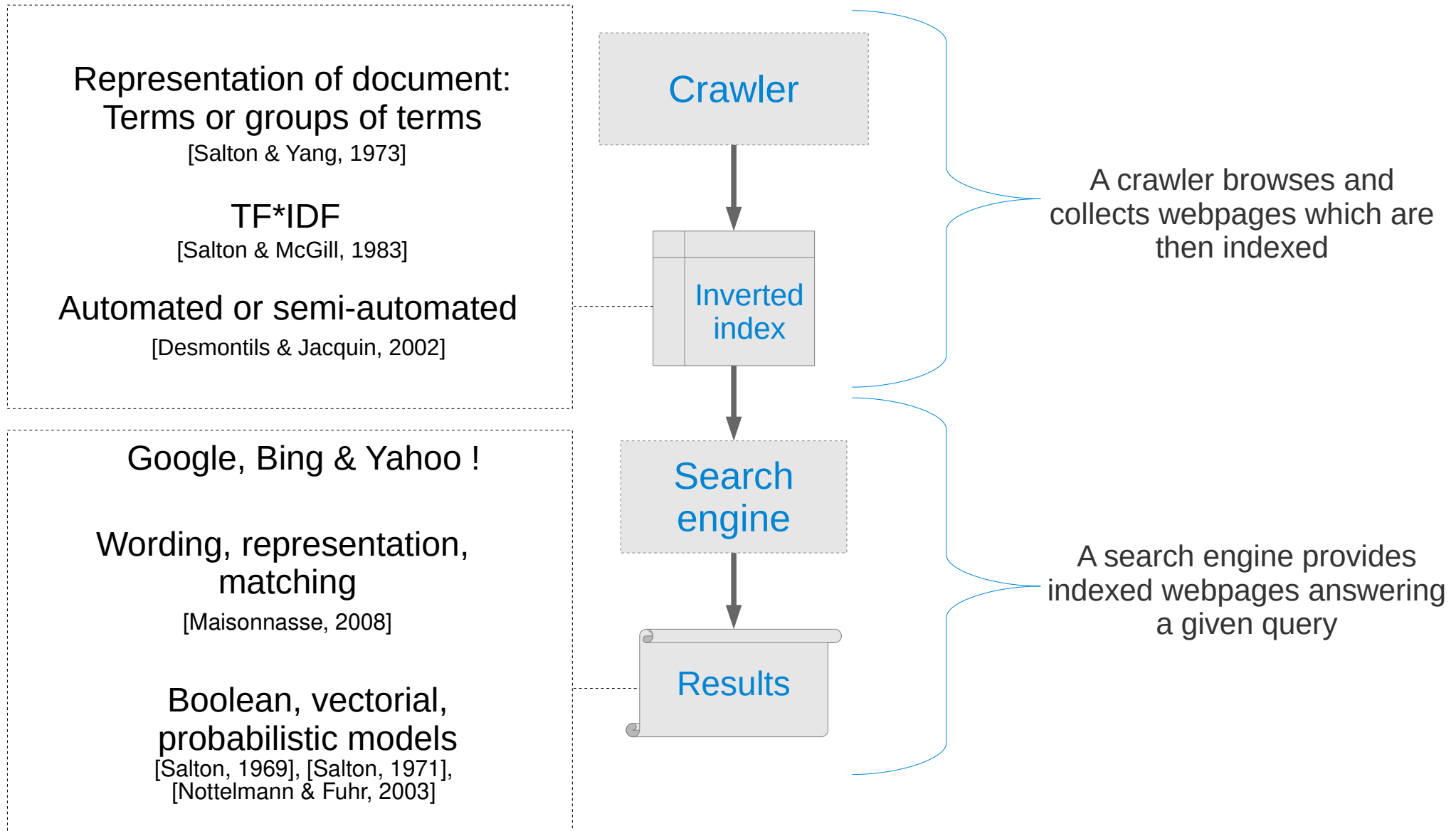
Web crawler



Inverted index

- **Definition:** index data structure used to store a set of documents or elements, mapping them from their content such as words or numbers.
- **Examples:**
 - "crepes" | {flour, egg, milk}
 - "genoise" | {egg, sugar, flour}
 - "caramel" | {sugar, butter}
 - "custard" | {egg, milk, sugar}
 - "flour" | {crepes, genoise}
 - "egg" | {crepes, genoise, custard}
 - "milk" | {crepes, custard}
 - "sugar" | {genoise, caramel, custard}
 - "butter" | {caramel}

Status of IR over the Web



Boolean model

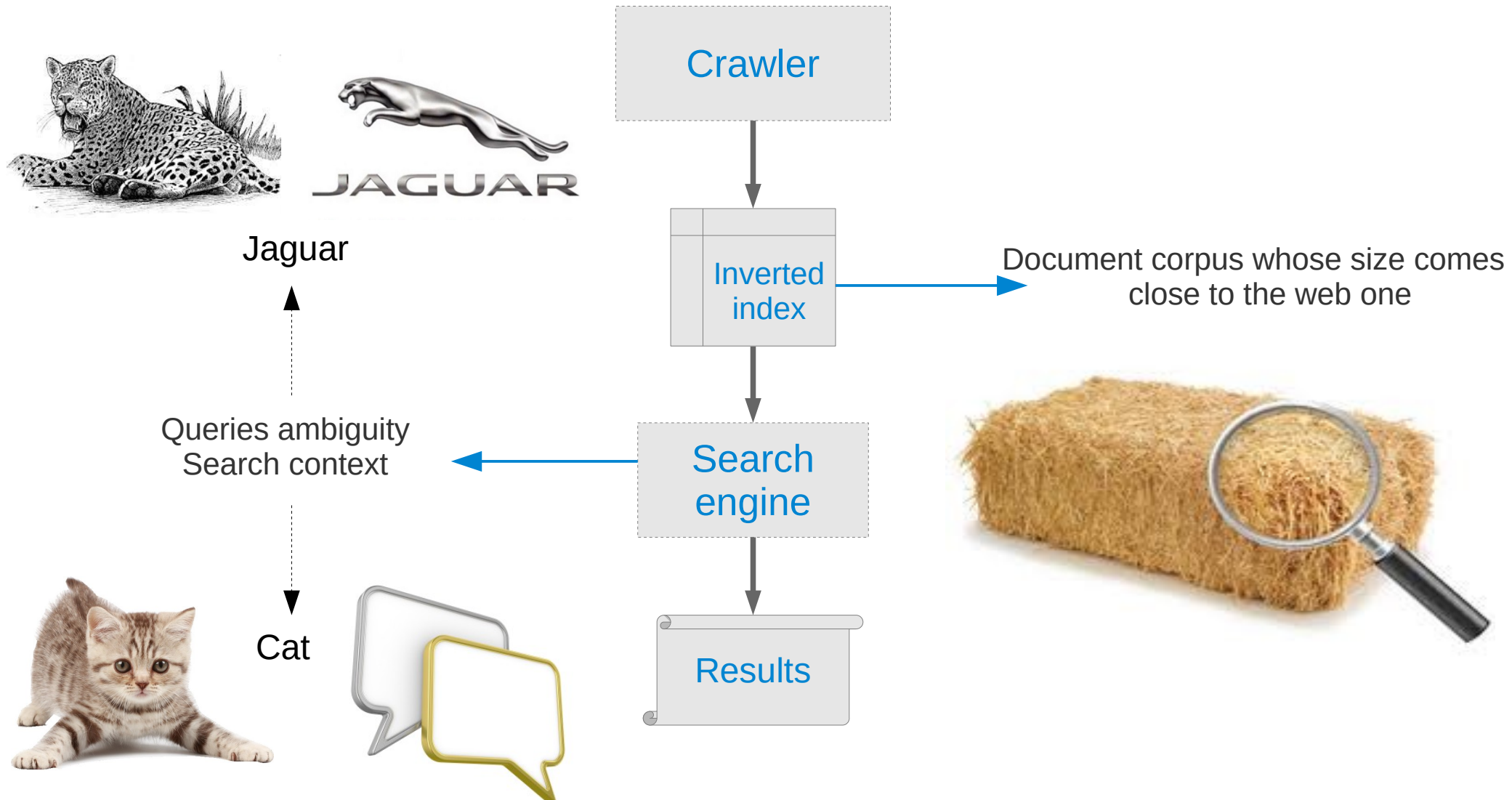
- First IR model, based upon set theory.
- A document is represented as a set of terms
 - Example: $d1(t1, t2, t5)$; $d2(t1, t3, t5, t6)$; $d3(t1, t2, t3, t4, t5)$
- Query = set of words and boolean operators (\wedge , \vee , \neg)
 - Example: $q = t1 \wedge (t2 \vee \neg t3)$
- **Exact** document/query matching relies on the presence or absence of query terms in the document
 - Example: $\text{Match}(q, d1)=1$; $\text{Match}(q, d2)=0$
- Drawbacks:
 - Selection of a document is a binary decision
 - Selected documents are not sorted
 - Query formulation is difficult for numerous users
 - Size problem: a lot of documents are returned

Vectorial model

- Proposed by Salton in the SMART system (Salton G., 1970).
- Working principle: documents and queries are represented as vectors in the space of document collection terms.
 - Document j : $d_j = (w_{1j}, w_{2j}, \dots, w_{Mj})$
 - Query: $q = (w_{1q}, w_{2q}, \dots, w_{Mq})$
 - With w_{ij} : weight of term t_i in document d_j (e.g.: tf*idf)
- A collection of n different documents and M different terms can be represented as a $n \times M$ matrix, and the query as a vector
- Relevance is a similarity measure between vectors (e.g.: cosine similarity)
- Pros:
 - Weighting enhances search results
 - The similarity measure allows to sort the documents by relevance
- Cons:
 - Do not take into account the word order (bag of words)
 - Vectorial representation necessitates term independence

Model which leads
to RAG !

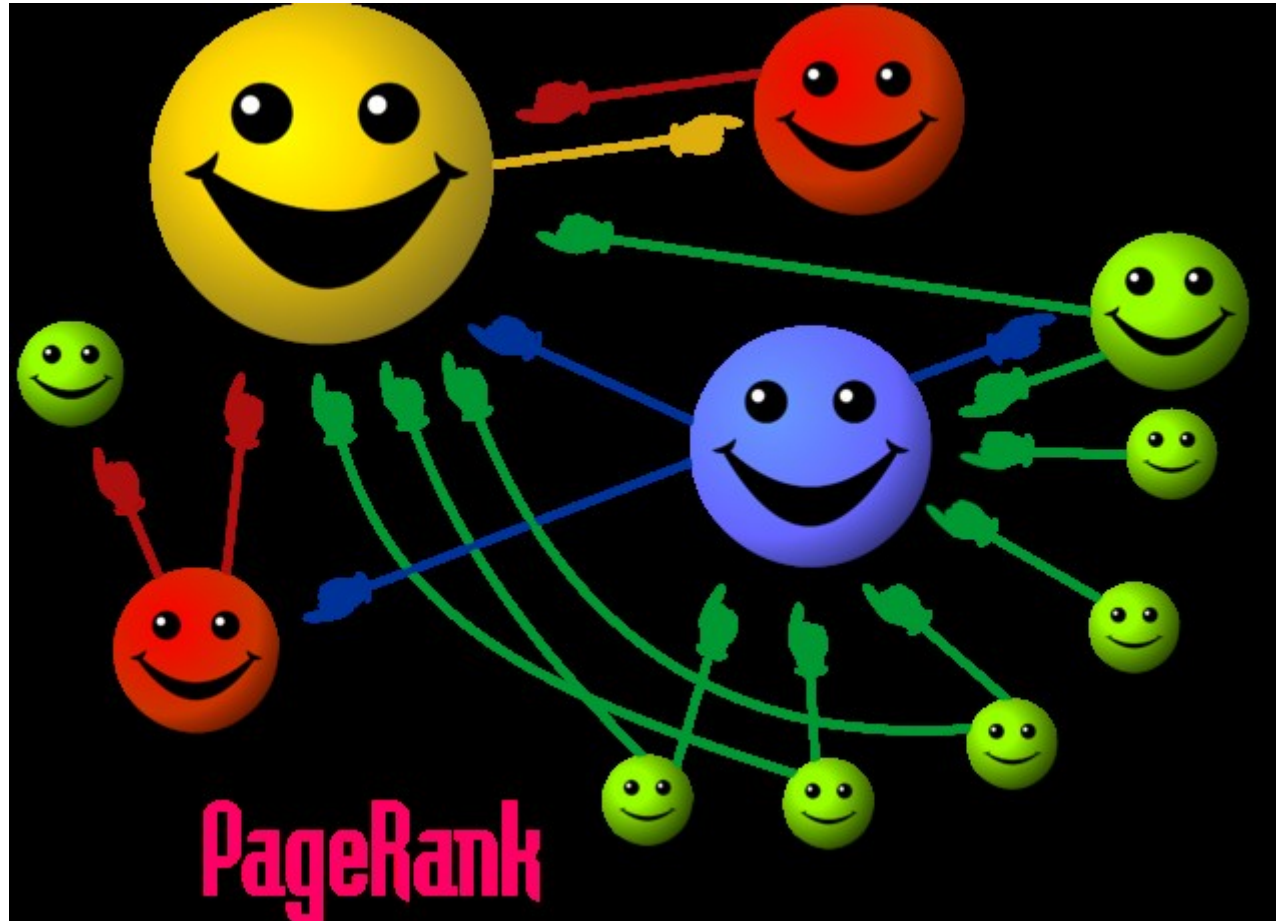
Limitations of IR on the web



Indexing

- Analysis (parsing) of pages/documents
 - Depends on the type of document
 - Conversion into **Hits** = weighting according to:
 - Number of appearances of each word
 - Position in the document
 - Relative font size and word letter case
- Only the surface web is indexed
 - Deep web: not-linked-to websites, restricted access content, content avoiding indexing, dynamic webpages
 - Deep web = 500 times bigger than surface web according to BrightPlanet, 2001

Google: PageRank (Page et al., 1998)



<http://en.wikipedia.org/wiki/File:PageRank-hi-res.png>

A webpage is important if linked to other important webpages

PageRank ~ probability to access a webpage by pure chance

PageRank formula

3 factors determine PageRank of a webpage P:

- Number of links pointing to P
- Number of links contained into webpages that have a link to P
- PageRank of pages containing a link to P
- PageRank is the summation of PageRanks of pages containing a link to P_i , weighted by the total number of output links:

$$r(P_i) = \sum_{P_j \in \mathcal{B}_{P_i}} \frac{r(P_j)}{|P_j|}$$

- Calculation is iterative

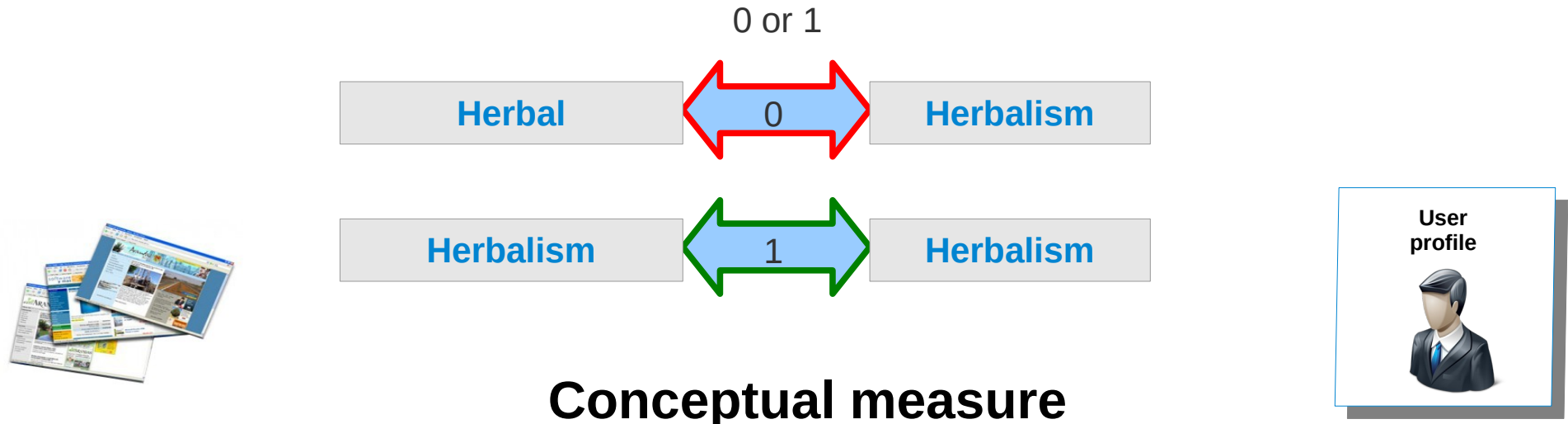
IR at Google's

Indexing documents

- Inverted index (word(s) → document(s))
- Relevance of words weighted by their type (title, anchor, URL, font...) and their position into the webpage
- Equivalent to $TF*IDF$
- Query/document similarity measure (dot product of the 2 vectors)
- Combination of IR score with PageRank in order to sort the results

Examples of relevance measure

Measure relative to terms (cosine similarity)



Basic HMI

- Minimalistic HMI: 1 text field, 1 result page (snippet) including extracts
- Advanced search is usually proposed

The screenshot shows a search engine interface with a search bar containing 'Rouen cathedral'. Below the search bar are tabs for 'All', 'Images', 'Videos', 'Maps', 'News', 'Shopping', 'My saves', and 'About search results'. The 'All' tab is selected. The main result is a Wikipedia entry titled 'Rouen Cathedral (Monet series) - Wikipedia' with the URL [https://en.wikipedia.org/wiki/Rouen_Cathedral_\(Monet_series\)](https://en.wikipedia.org/wiki/Rouen_Cathedral_(Monet_series)). Below the title are tabs for 'Overview', 'Date', 'Painting Light', 'Technique', 'Gallery', and 'P1'. The 'Overview' tab is selected. The text under the 'Overview' tab reads: 'The Rouen Cathedral series was painted in the 1890s by French impressionist Claude Monet. The paintings in the series each capture the façade of the Rouen Cathedral at different times of the day and year and reflect changes in its appearance under different lighting conditions.' Below this text is a link 'See more on en.wikipedia.org · Text under CC-BY-SA license'. Further down are two rows of metadata: 'Artist: Claude Monet' and 'Location: Musée d'Orsay, Paris, France' in the first row, and 'Dimensions: 107 cm × 73.5 cm (42 in × 28.9 in)' and 'Year: 1894' in the second row. To the right of the main result is a section titled 'Related Searches for rouen cathedral' with eight search suggestions: 'notre dame cathedral rouen', 'rouen cathedral monet', 'rouen cathedral opening times', 'rouen cathedral music', 'rouen cathedral mass times', 'monet rouen cathedral series', 'rouen cathedral claude monet', and 'notre dame rouen'. At the bottom of the page is another result snippet for 'Rouen Cathedral - Wikipedia' with the URL https://en.wikipedia.org/wiki/Rouen_Cathedral and a tab for 'Overview'.

Rouen cathedral

All Images Videos Maps News Shopping | My saves About search results ⓘ

Rouen Cathedral (Monet series) - Wikipedia
[https://en.wikipedia.org/wiki/Rouen_Cathedral_\(Monet_series\)](https://en.wikipedia.org/wiki/Rouen_Cathedral_(Monet_series))

< Overview Date Painting Light Technique Gallery P1 >

The Rouen Cathedral series was painted in the 1890s by French impressionist Claude Monet. The paintings in the series each capture the façade of the Rouen Cathedral at different times of the day and year and reflect changes in its appearance under different lighting conditions.

[See more on en.wikipedia.org](#) · Text under CC-BY-SA license

Artist: [Claude Monet](#) Location: [Musée d'Orsay, Paris, France](#)

Dimensions: 107 cm × 73.5 cm (42 in × 28.9 in) Year: 1894

Related Searches for rouen cathedral

[notre dame cathedral rouen](#) [rouen cathedral monet](#)

[rouen cathedral opening times](#) [rouen cathedral music](#)

[rouen cathedral mass times](#) [monet rouen cathedral series](#)

[rouen cathedral claude monet](#) [notre dame rouen](#)

Rouen Cathedral - Wikipedia
https://en.wikipedia.org/wiki/Rouen_Cathedral

Overview

HMI: input help

- Spell check
- Automatic completion proposing the most searched queries



HMI: result analysis help

- Preview of results
- Sponsored links
- Terms and directories
- Language choice

The screenshot shows the EXALEAD search engine interface. At the top, there's a navigation bar with links for Web, Images, Videos, Wikipedia, and More. The search bar contains the text 'cathédrale de rouen' and a 'Search' button. Below the search bar, related searches are listed: 'Hotel Rouen', 'Rouen Cathedral', and 'Rouen Centre'. The main results area shows a list of search results. The first result is from Panoramio, showing a photo of the Cathédrale de Rouen. The second result is from Wikimedia Commons, showing a category page for Cathédrale Notre-Dame de Rouen. The third result is from TripAdvisor, showing hotel reviews for Hotel de la Cathedrale in Rouen. On the right side, there's a sidebar with filters for Site type (Blog, Forum), Filetype (pdf, rtf, swf, text, word), Related terms (Beaux Art, Champ elysees, Cities Center, Greater Britain, Hotel Vieux), and Languages (English (81%), French (14.8%), etc.).

Web Images Videos Wikipedia More ►

EXALEAD

cathédrale de rouen

Search

Advanced Search

Home > Web results 1-10 of 268,111 for **cathédrale de rouen**, Page 1 - Next page

Did you mean: **cathedral de rouen** ?

Panoramio - Photo of La Cathédrale de Rouen et l'Arc en Ciel.
Photo-sharing community. Discover the world through photos.
www.panoramio.com/photo/11999820
Cached - Bookmark

Category: Cathédrale Notre-Dame de Rouen - Wikimedia Commons
Category: **Cathédrale Notre-Dame de Rouen** From Wikimedia [...] Media in category "**Cathédrale Notre-Dame de Rouen** [...] Place de la **cathédrale (Rouen)** Monuments...
commons.wikimedia.org/wiki/Category:Cathédrale_Notre-Dame_de_Rouen
07 Sep 2013 - Cached - Bookmark

Hotel de la Cathedrale (Rouen, France) - Hotel Reviews - TripAdvisor
Hotel de la **Cathedrale, Rouen**: See 148 traveler reviews, 85 candid photos, and great deals for Hotel de la **Cathedrale**, ranked #17 of 52 hotels in
www.tripadvisor.com/Hotel_Review-g187191-d219860-Reviews-Hotel_de_la_Cathedrale-Rouen_Sein...

Site type:
» Blog
» Forum

Filetype:
» pdf
» rtf
» swf
» text
» word

Related terms:
» Beaux Art
» Champ elysees
» Cities Center
» Greater Britain
» Hotel Vieux

Languages :

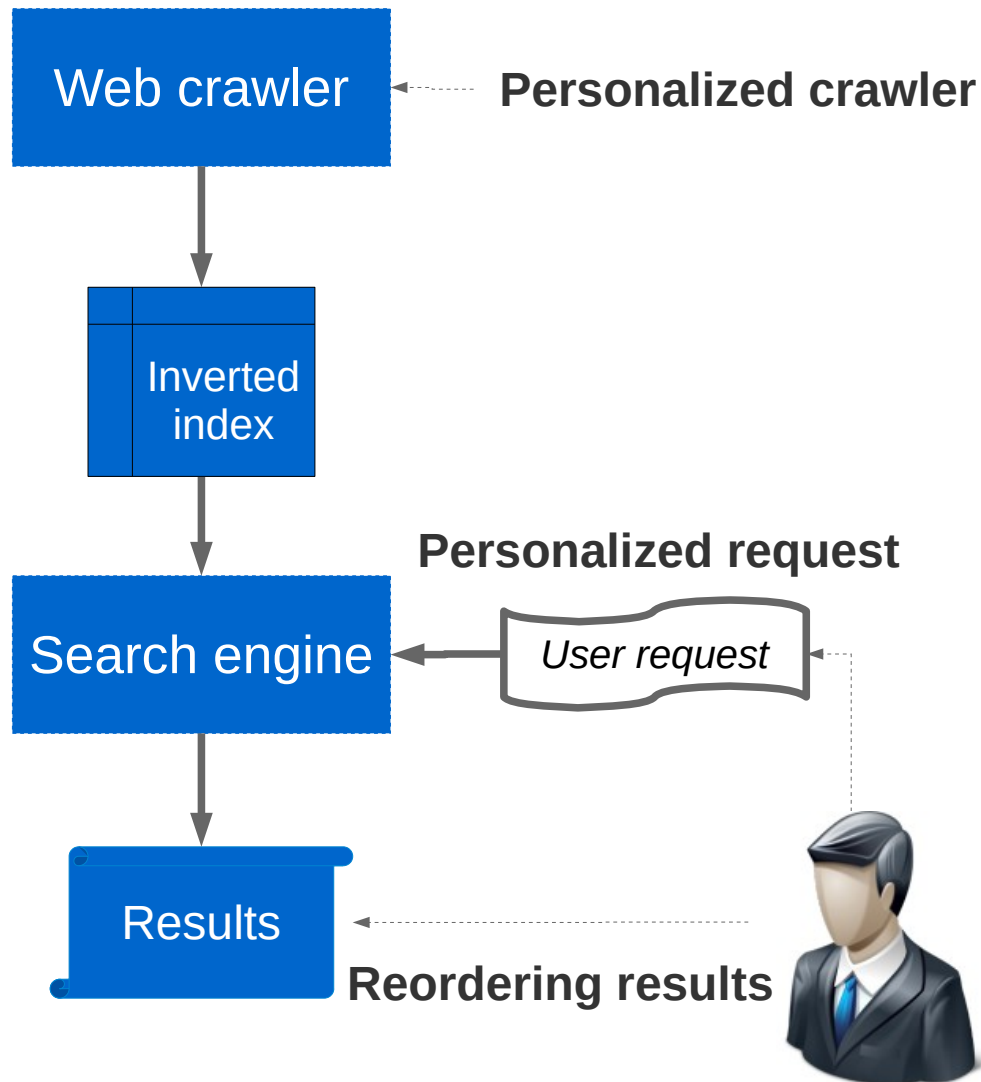
English (81%)
French (14.8%)

Drawbacks of classic approaches

- PageRank problems
 - Google: 82% of queries in France
 - A webpage is well referenced if popular ; a webpage is popular if well referenced
 - Difficult appearance of new pages
 - SEO is very lucrative
- Content quality is not taken into account
- Ambiguities? Synonymies?
- Search context?
- Generalization (concepts) ?

⇒ **Personalization!**

Personalization: founding approaches



Statistical model

C. Aggarwal, F. Al-Garawi, and P. Yu. Intelligent crawling on the world wide web with arbitrary predicates. 2001.

Reinforcement learning

S. Chakrabarti, K. Punera, and M. Subramanyam. Accelerated focused crawling through online relevance feedback. 2002.

Multi-agent : genetic model

F. Menczer and R. Belew. Adaptive retrieval agents : Internalizing local context and scaling up to the web. 2000.

Multi-agent : biological model

F. Gasparetti and A. Micarelli. Swarm intelligence : Agents for adaptive web search. In ECAI, 2004.

System – Human interaction in IR

- Example: R="animal", D1="horse"
 - D1 will not be returned as it does not contain "animal"
 - How to return relevant documents even if they do not contain any of the query terms?

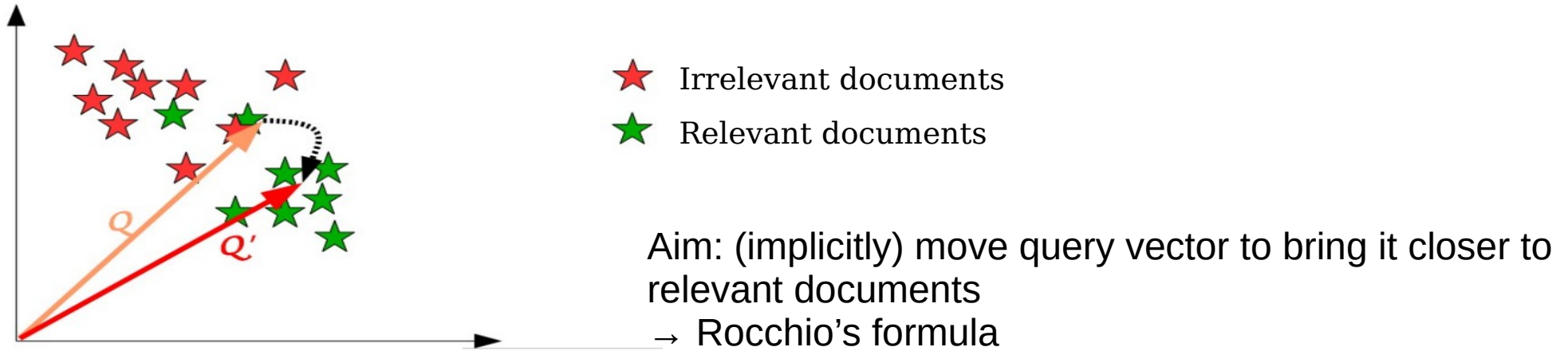
Methods improving recall:

- Query expansion
- Relevance feedback
- Query disambiguation: concepts, detection of synonyms, query validation through H-M "dialogue"...

Query expansion

- Query modification using external resources
 - Dictionaries, thesauri, ontologies...
 - Increase recall
 - Precision will lower consequently
- Examples:
 - Hospital → medical (thematic)
 - Interest rate → interesting (grammatical)
 - Shiny → bright (synonymy)
 - Horse → animal (generalization)
 - Animal → horse, dog, ... (specialization)

Relevance feedback



$$\vec{Q}' = \alpha \vec{Q} + \beta \vec{P} + \gamma \vec{N} \vec{P}$$

→ Mean of irrelevant documents vectors

→ Negative value (e.g.: -0.25)

→ Mean of relevant documents vectors

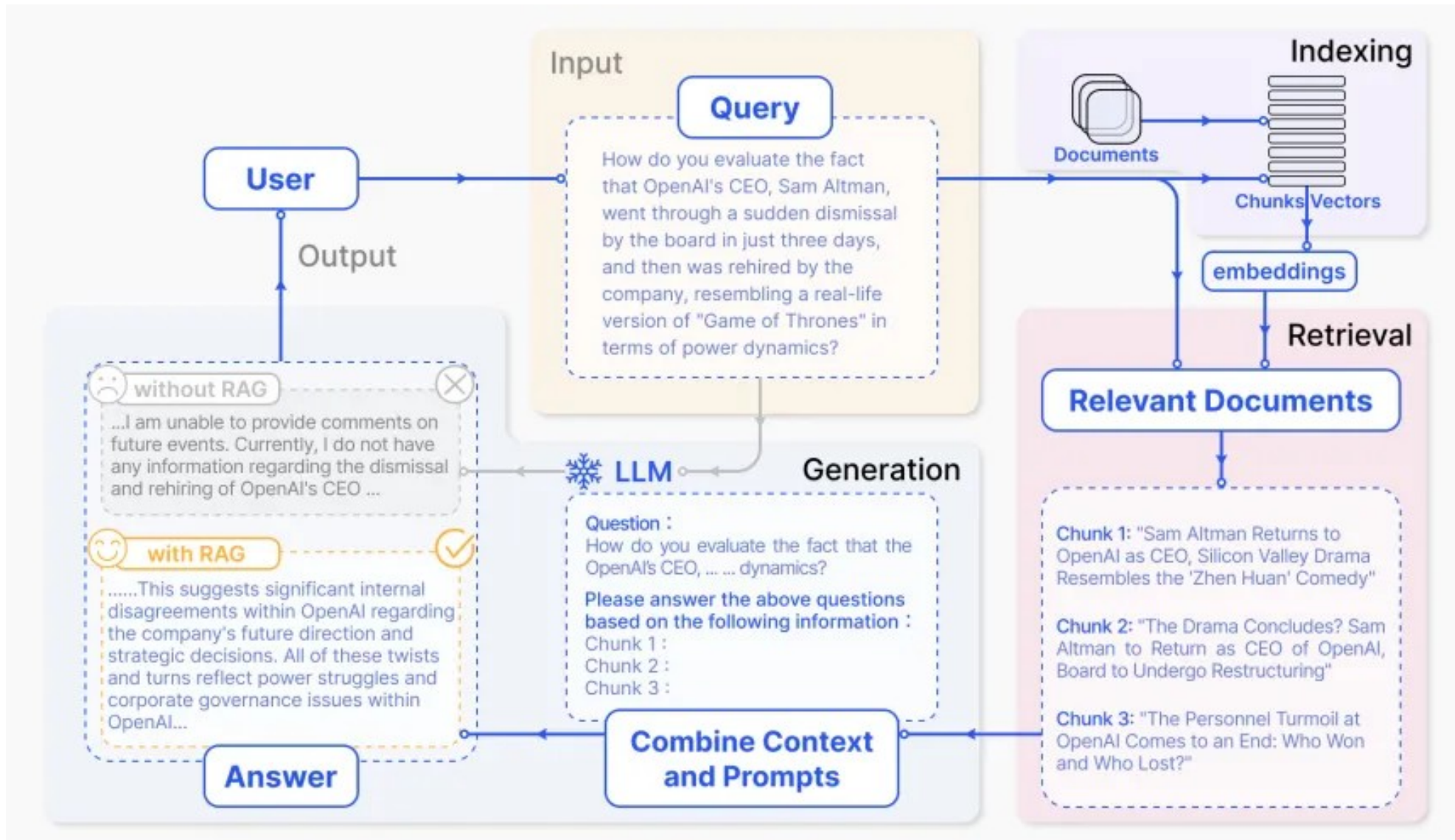
→ Positive value (e.g.: 0.5)

→ Initial query vector

→ Positive value greater than the others

→ New query vector

Personalised RAG?



Hugely inspired by

Lectures

- Xavier Tannier's lectures
http://perso.limsi.fr/amax/enseignement/iri/M2PRO_IRI_6_RIWeb.pdf
- <http://math.univ-lyon1.fr/homes-www/malbos/lib/exe/fetch.php?media=ens:algapp12:algapiichapiiisec9.pdf>
- <http://www.iro.umontreal.ca/~nie/IFT6255/Introduction.html>
- https://www.irit.fr/~Mohand.Boughanem/Enseignements_RI.php