

Temu-Kembali Informasi 2018

02: Arsitektur Search Engine

Versi Ringkas ++

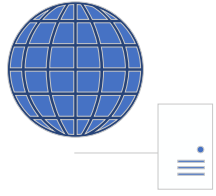
Arsitektur Software

- Arsitektur Software merujuk ke struktur tingkat tinggi dari suatu sistem perangkat lunak.
- Struktur ini diperlukan untuk menjelaskan tentang sistem perangkat lunak.
- Setiap struktur terdiri dari elemen perangkat lunak, hubungan di antara mereka, dan properti dari elemen dan relasi tersebut.
- [Wikipedia]

Contoh 1:

Arsitektur Search Engine







Proses Indexing



Data Storage



Proses Pencarian (*Search*)



Akuisisi



Konversi ke plain text dan unified encoding

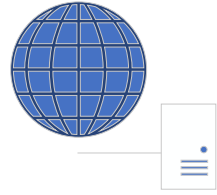
Proses Indexing



Data Storage



Proses Pencarian (*Search*)



Akuisisi

Konversi ke plain text dan unified encoding

Transformasi

Index terms, fitur, klasifikasi, meta data

d
t₁
t₂
t₃
...
f₁, f₂, f₃, ...
c₁ not spam
c₂ sports
...
o₁ 10 inlinks
...

Proses Indexing



Data Storage



Proses Pencarian (*Search*)



Akuisisi

Konversi ke plain text dan unified encoding

Transformasi

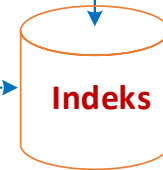
Index terms, fitur, klasifikasi, meta data

d
 $t_1 \begin{pmatrix} 0.1 \\ 0.3 \\ 0.2 \\ \vdots \end{pmatrix}$
 t_2
 t_3
 \vdots
 f_1, f_2, f_3, \dots
 c_1 not spam
 c_2 sports
 \vdots
 o_1 10 inlinks
 \vdots

Indexing

Statistika, Pembobotan

Proses Indexing

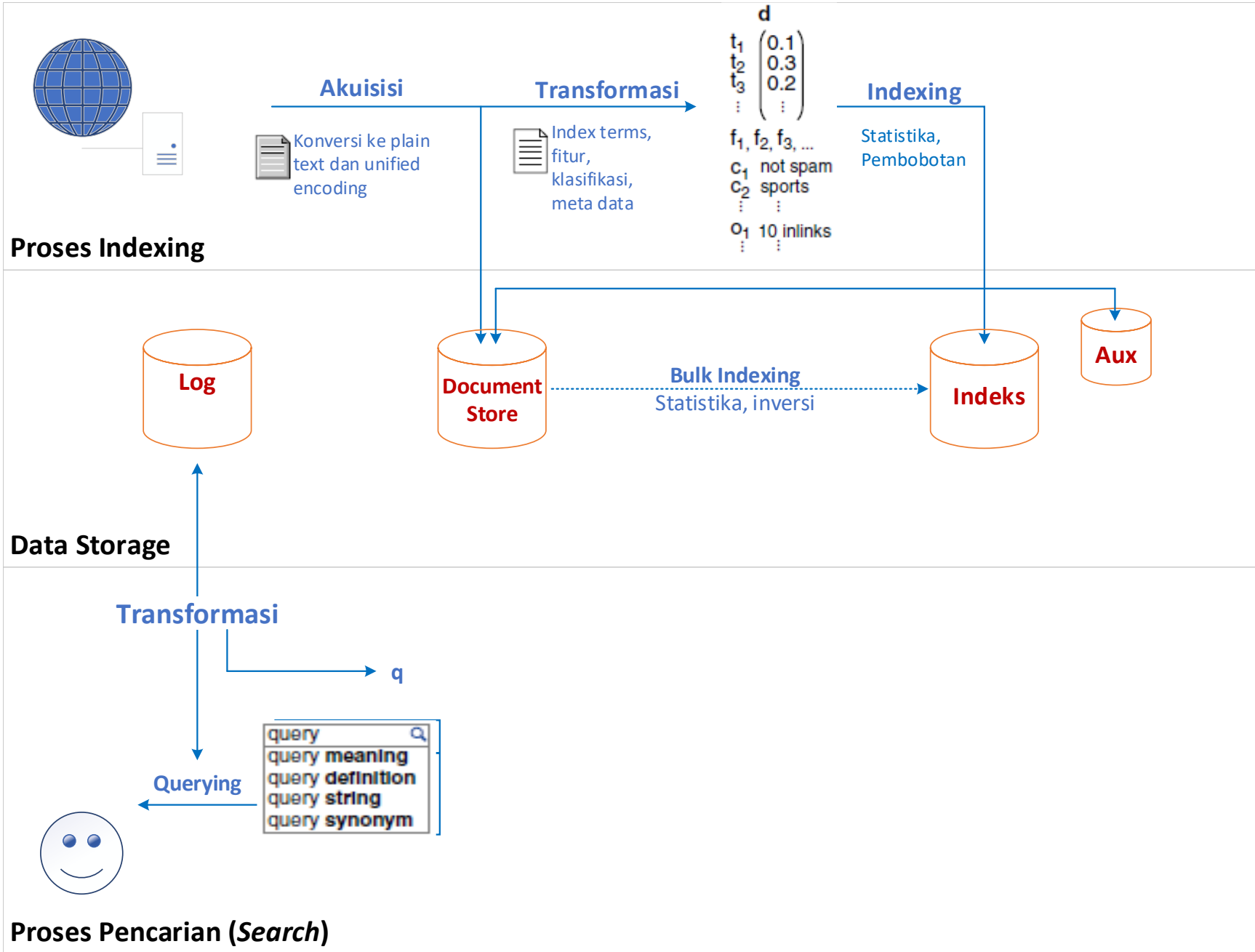


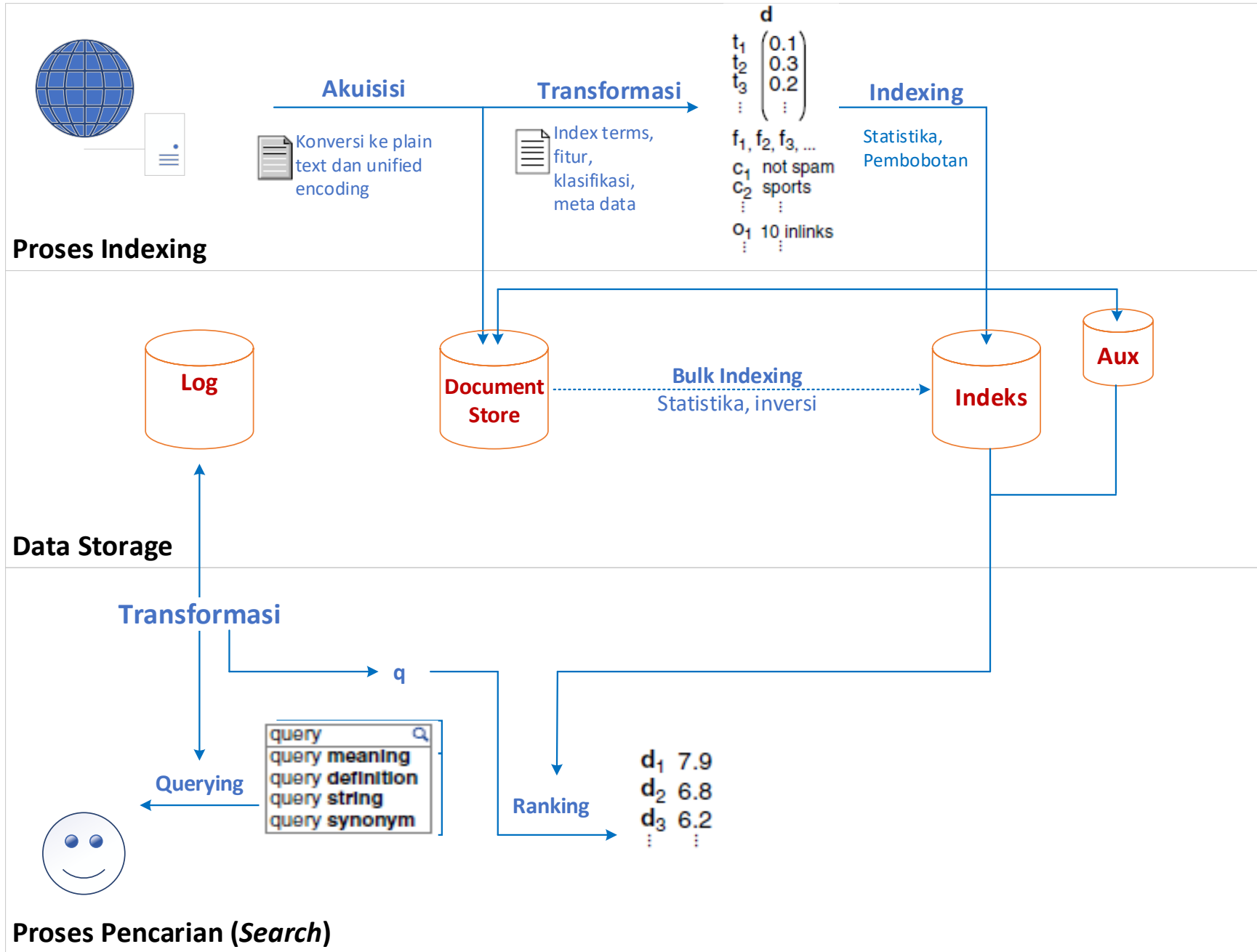
Bulk Indexing
Statistika, inversi

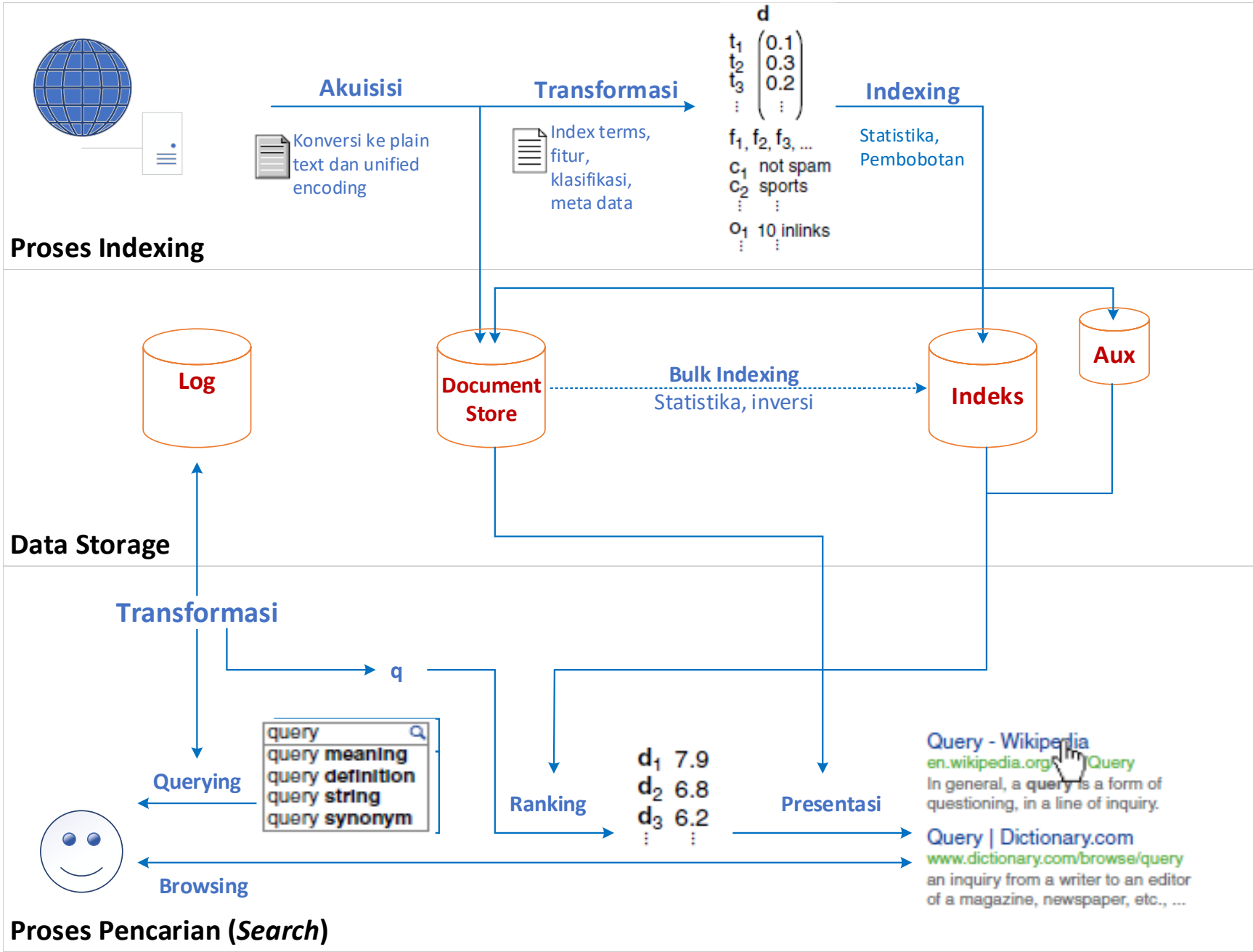
Data Storage

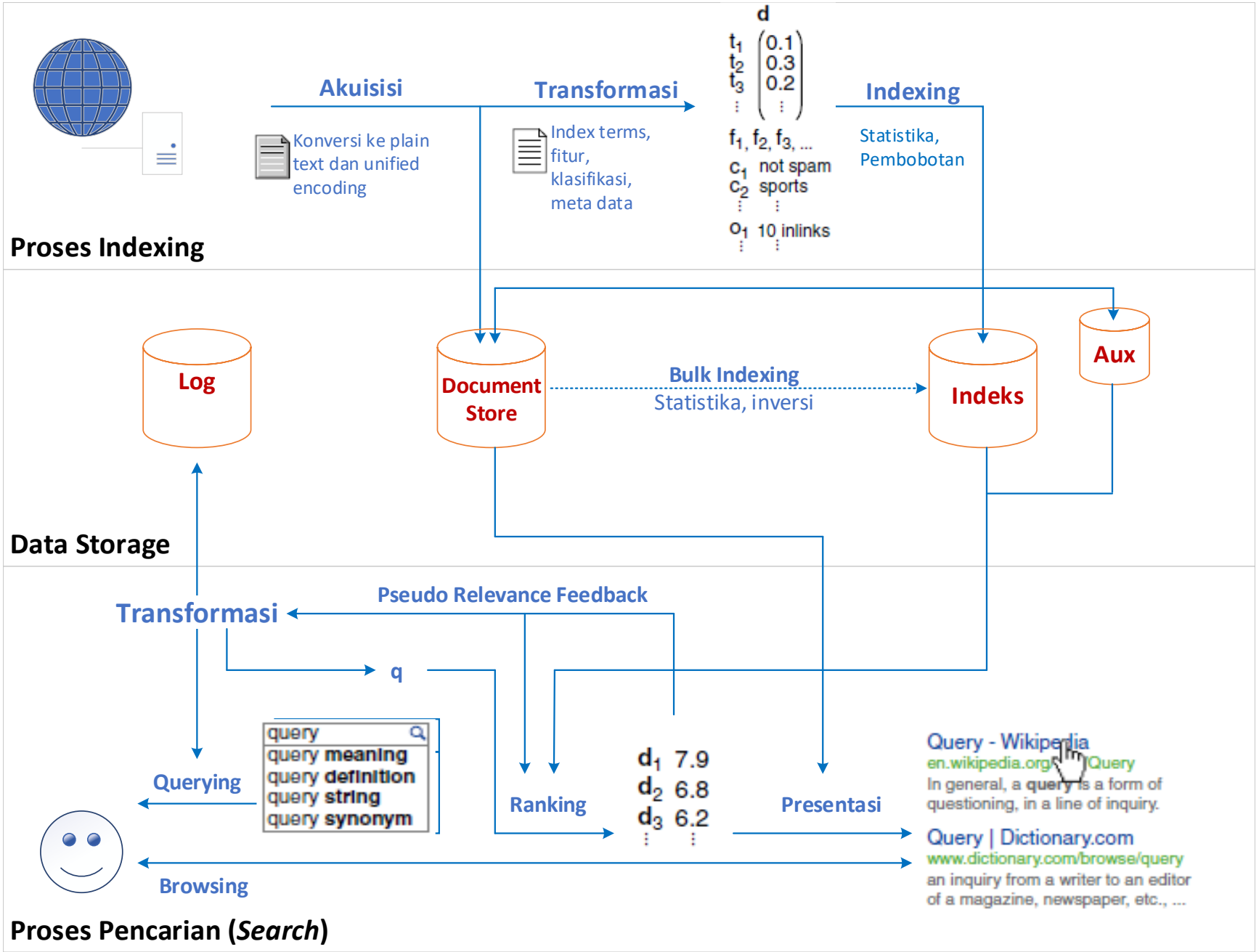


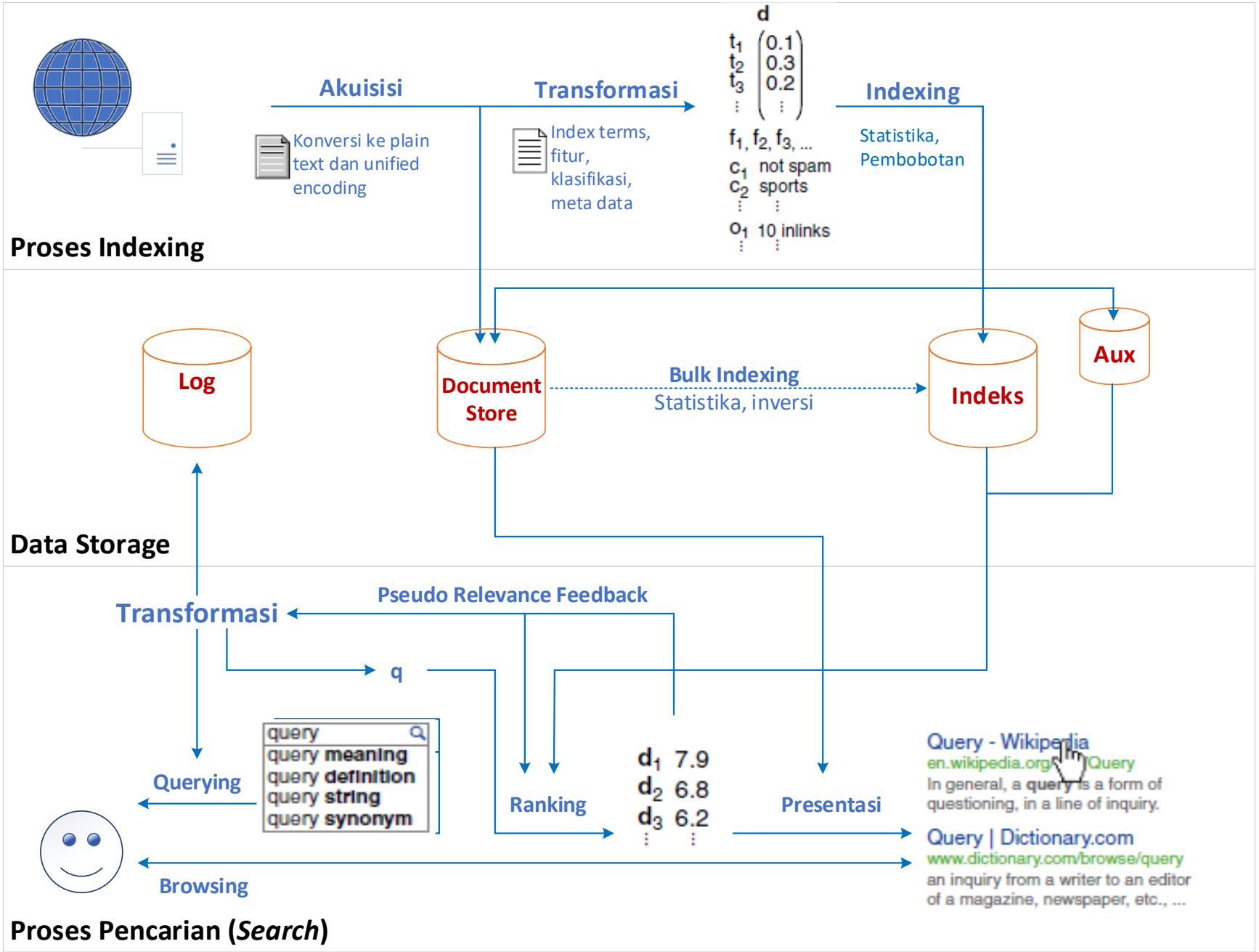
Proses Pencarian (Search)

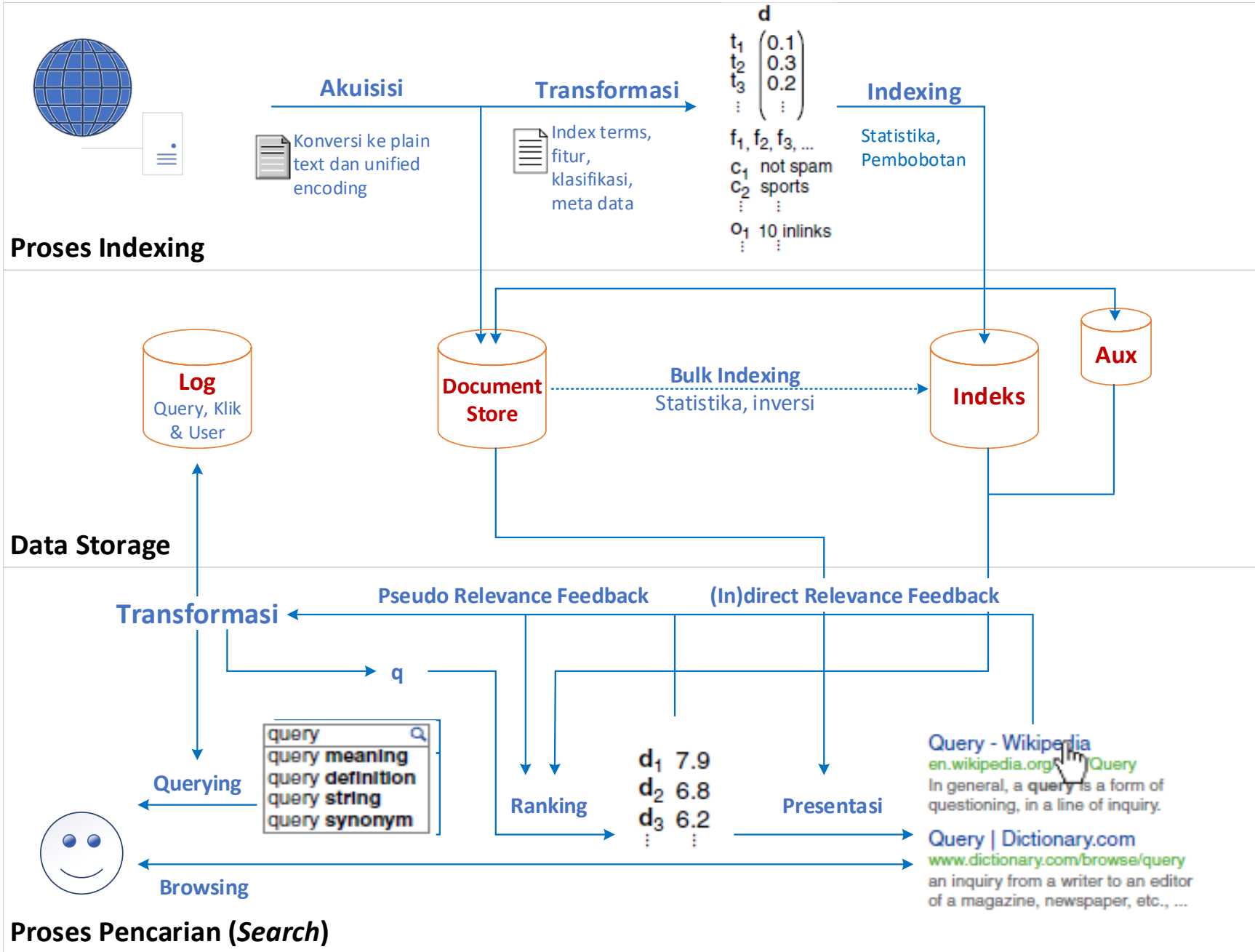














Akuisisi

Konversi ke plain text dan unified encoding

Transformasi

Index terms, fitur, klasifikasi, meta data

Indexing

Statistika, Pembobotan

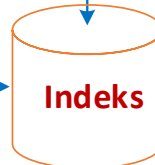
$$d \begin{pmatrix} 0.1 \\ 0.3 \\ 0.2 \\ \vdots \end{pmatrix}$$

t_1
 t_2
 t_3
 \vdots
 f_1, f_2, f_3, \dots
 c_1 not spam
 c_2 sports
 \vdots
 \vdots
 o_1 10 inlinks
 \vdots

Proses Indexing



Bulk Indexing
Statistika, inversi



Data Storage

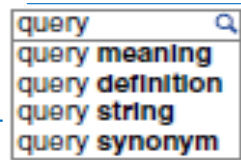
Transformasi

Pseudo Relevance Feedback

(In)direct Relevance Feedback

q

Querying



Ranking

$$d_1 \ 7.9$$

$$d_2 \ 6.8$$

$$d_3 \ 6.2$$

$$\vdots \ \vdots$$

Presentasi

Query - Wikipedia
en.wikipedia.org/wiki/Query
 In general, a **query** is a form of questioning, in a line of inquiry.

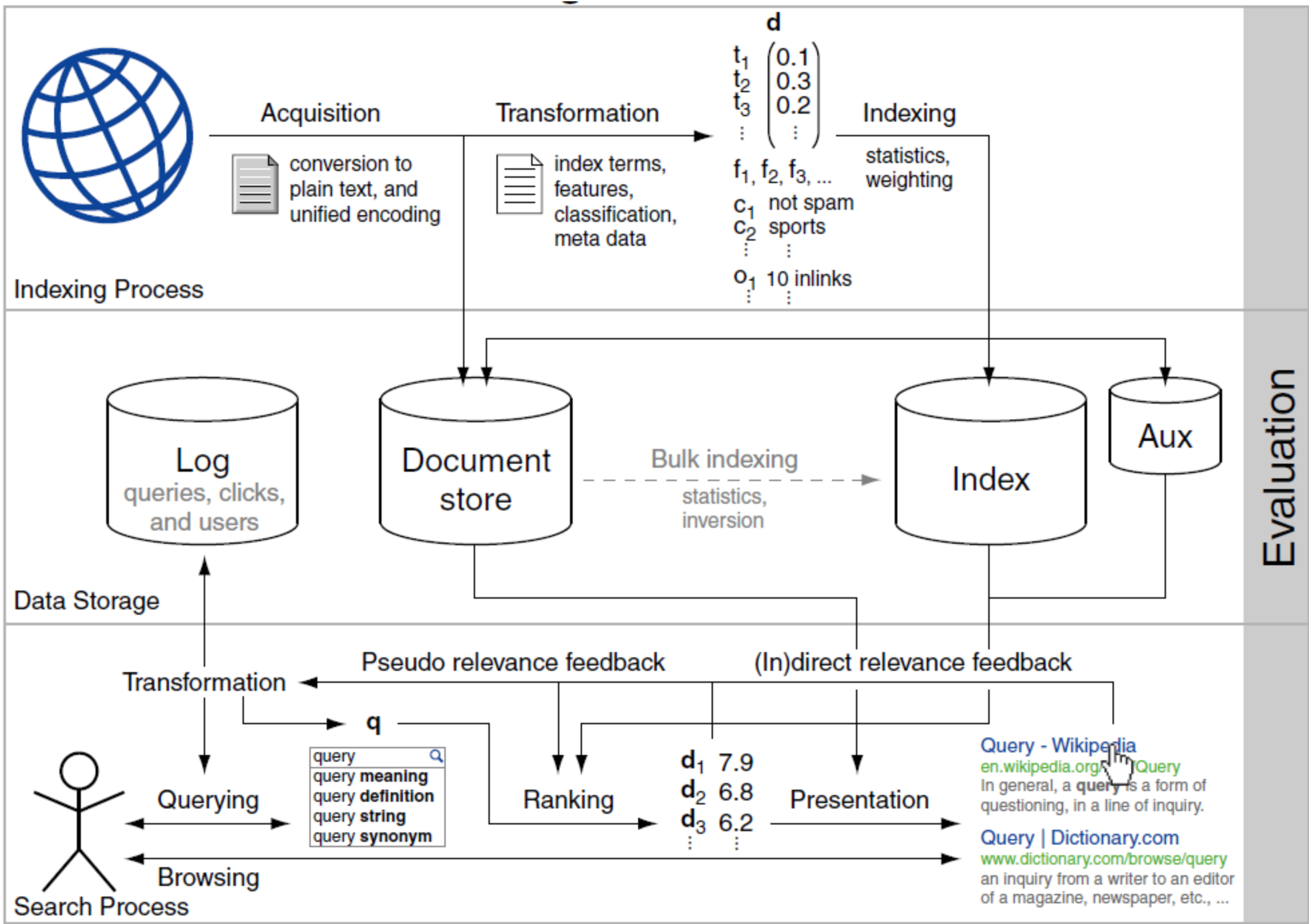
Query | Dictionary.com
www.dictionary.com/browse/query
 an inquiry from a writer to an editor of a magazine, newspaper, etc., ...



Browsing

Proses Pencarian (Search)

EVALUASI

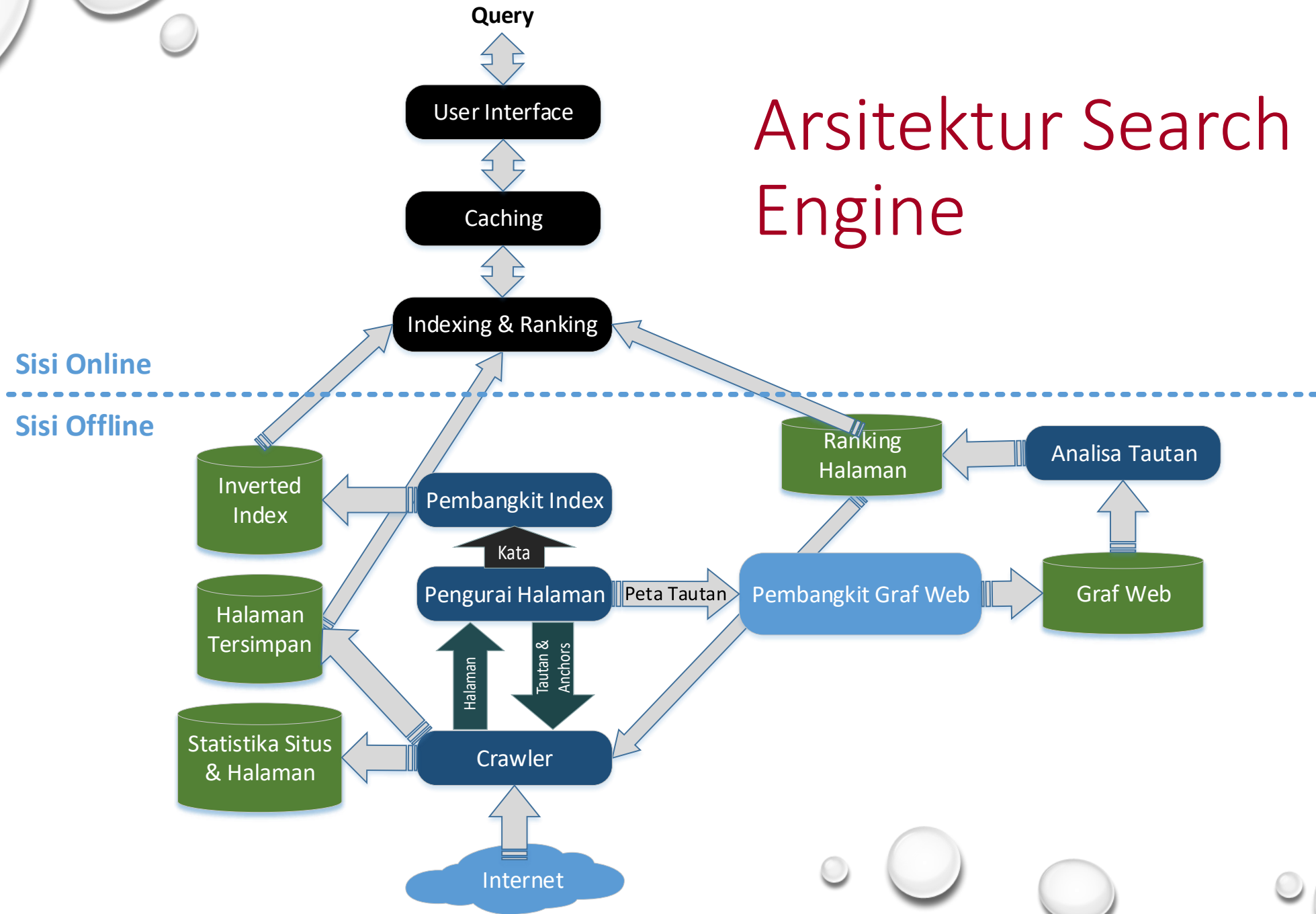


Evaluation

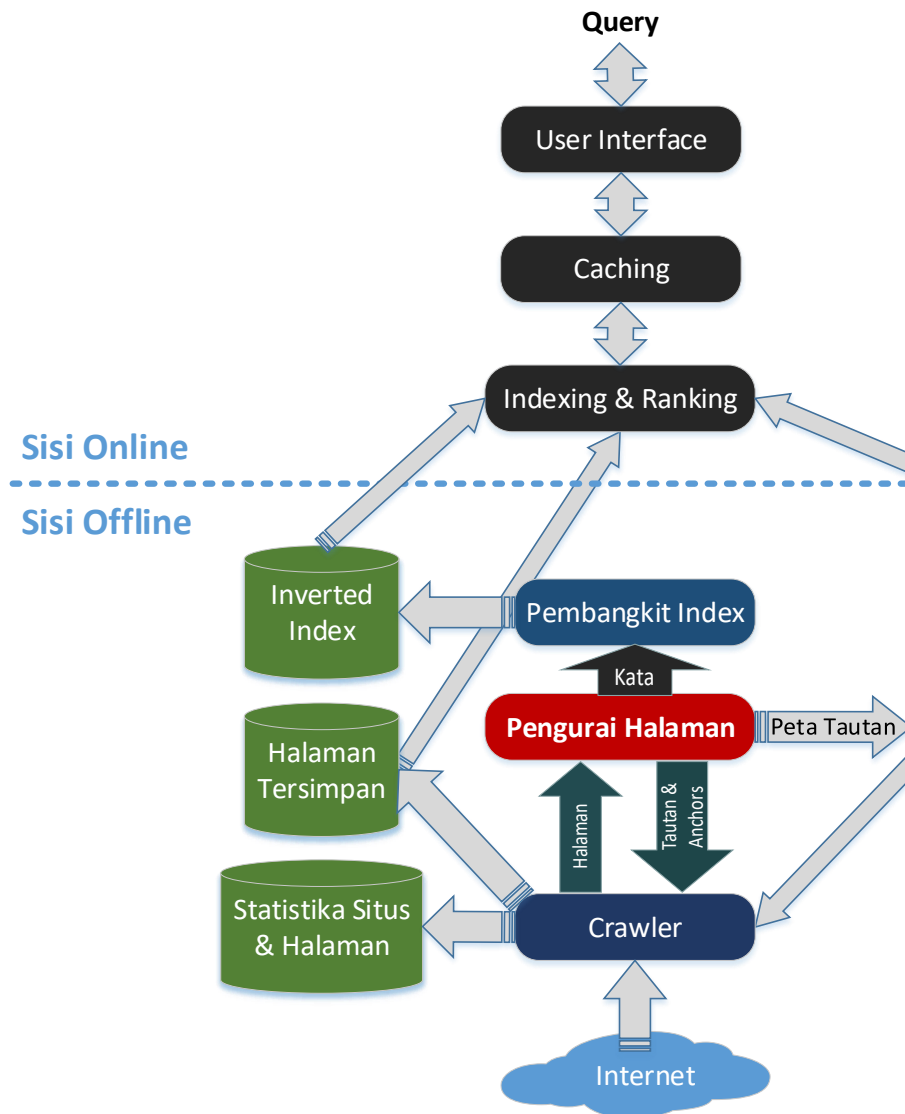
Contoh 2:

Arsitektur Search Engine dari Microsoft Research

Arsitektur Search Engine



Arsitektur: Page Parser



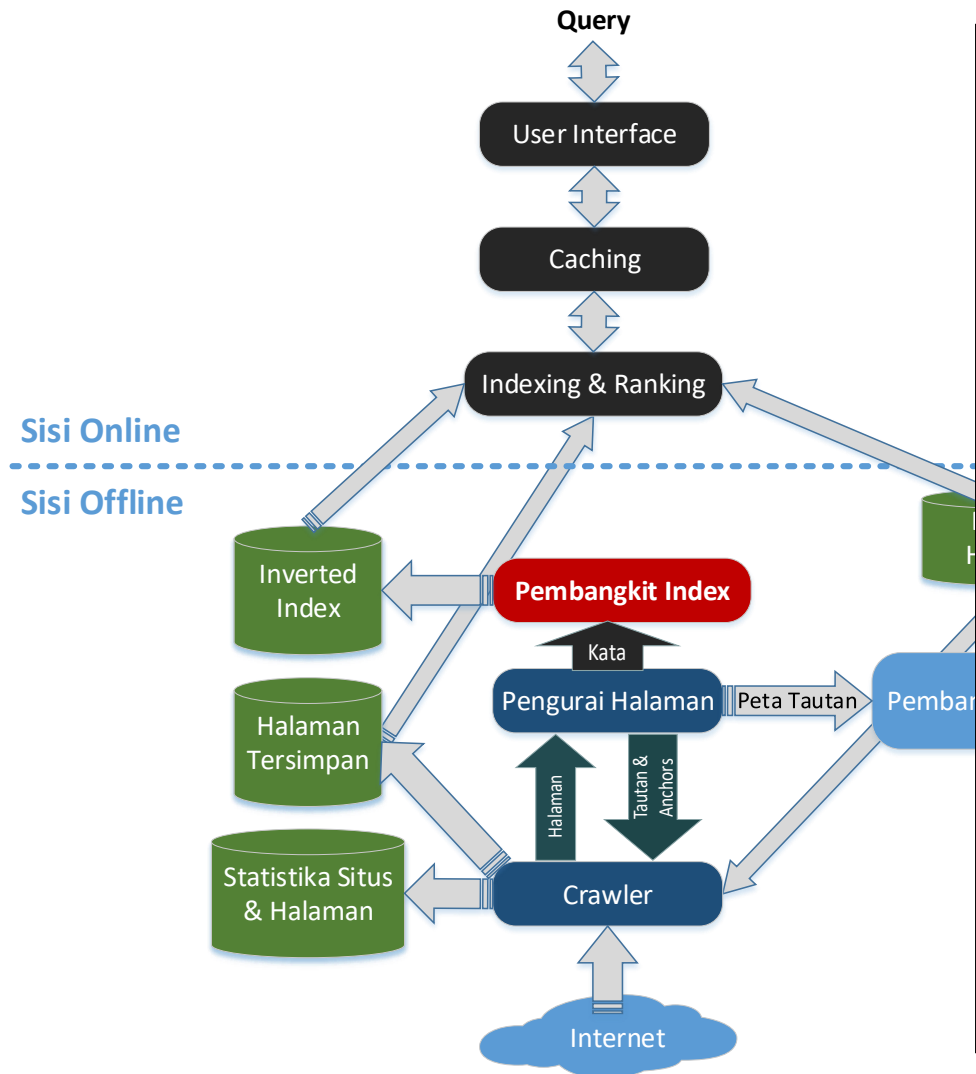
➤ Fungsi

- Mengekstrak aliran data untuk indexing
 - a. Title: kata-kata dalam <title>...</title>
 - b. URL
 - c. Body
 - Teks Anchor
 - Teks Plain
 - H1...6
 - Bold, Italic, etc
 - Large, Medium, Small
- Membangun peta link parsial
- Mengirim hyperlink yang ditemukan ke crawler

➤ Masalah Inti

- Fitur apa yang akan diekstrak?

Arsitektur: Index Builder



➤ Fungsi

- Membangun inverted index berdasarkan pada data halaman yang telah diparse

TermID	DocNum
	DocID HitNum Hit Hit Hit ...
	DocID HitNum Hit Hit Hit ...

➤ Masalah Inti

- Efisiensi vs. memory terbatas & terdistribusi

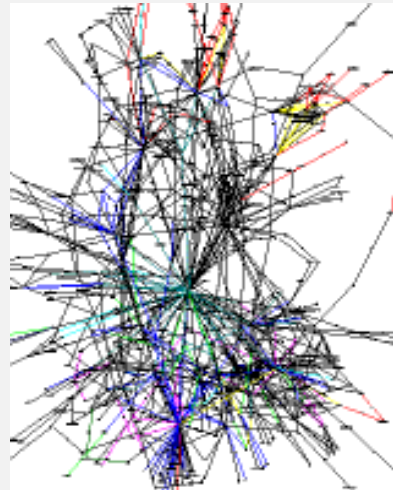
➤ Solusi

- Indexing terdistribusi
- Partisi berdasarkan dokumen, bukan partisi berdasarkan term

Arsitektur: Link Analysis

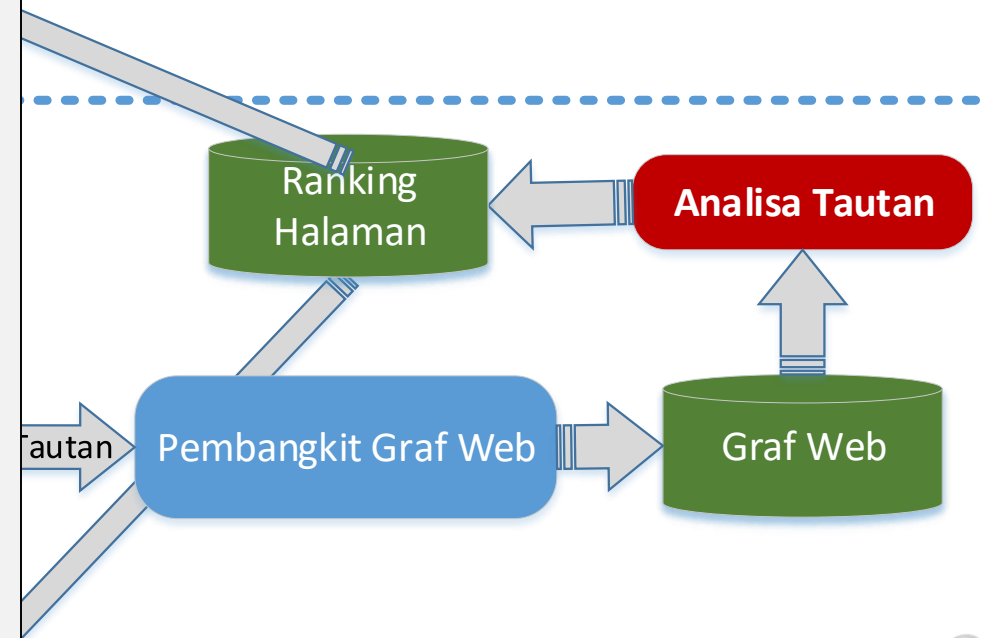
➤ Fungsi

- Mengukur kualitas atau otoritas dari suatu halaman berdasarkan pada graf link

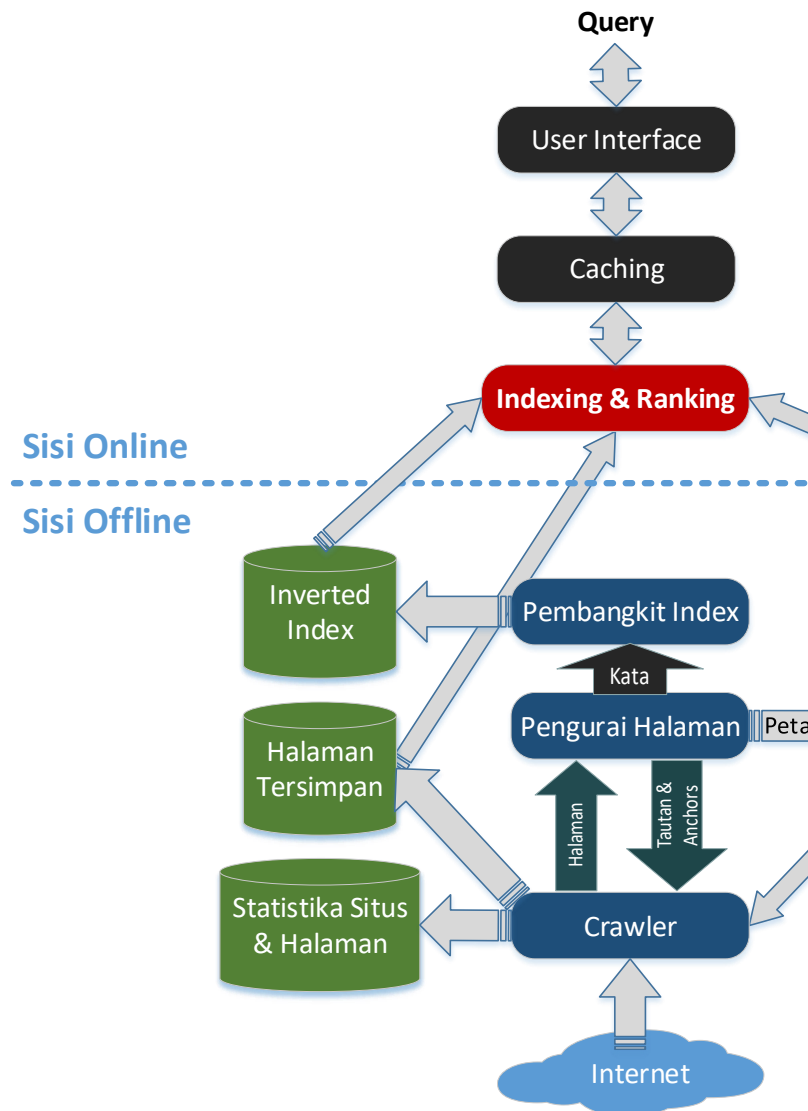


➤ Masalah Inti

- Algoritma yang efisien pada graf raksasa
- Link-spam?
- Apakah hanya link analysis cara untuk menentukan kualitas dari halaman?

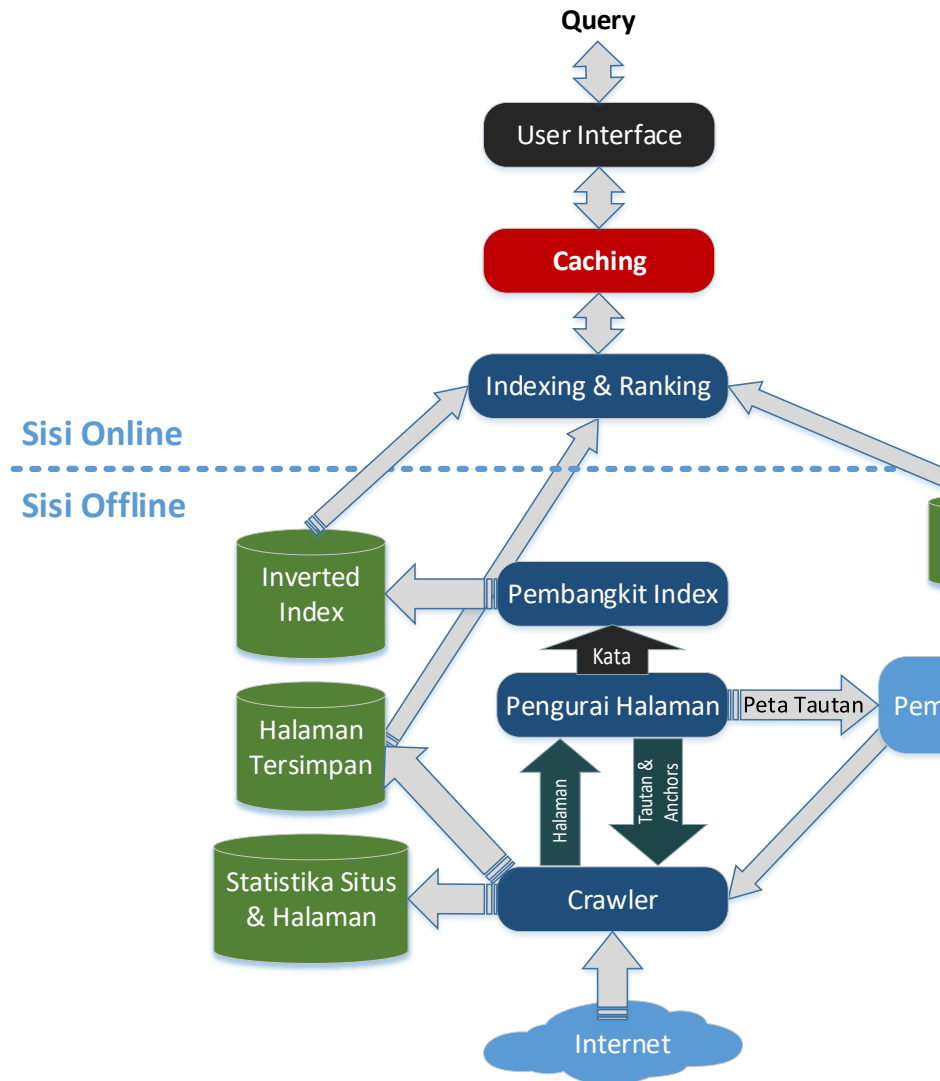


Arsitektur: Indexing & Ranking



- Masalah utama dalam komunitas IR dan telah dikaji puluhan tahun
- Fungsi
 - Indexing: dengan cepat menemukan halaman yang mengandung term query
 - Ranking: mengurutkan halaman sesuai dengan relevansi terhadap query
- Masalah Inti
 - Kinerja: inverted list untuk suatu term hot mungkin ratusan megabyte.
 - Akurasi: fungsi ranking dengan ratusan parameter:
 - Teks Anchor
 - Ranking halaman
 - Term proximity
 - TF*IDF
 - ...
- Solusi
 - Kinerja: Top-K query & index pruning
 - Akurasi: Tuning atau learning?

Arsitektur: Caching



➤ Fungsi

- Men-cache hasil dari query yang sering untuk menjawab ribuan query per detik dengan waktu respon interaktif

➤ Masalah Inti

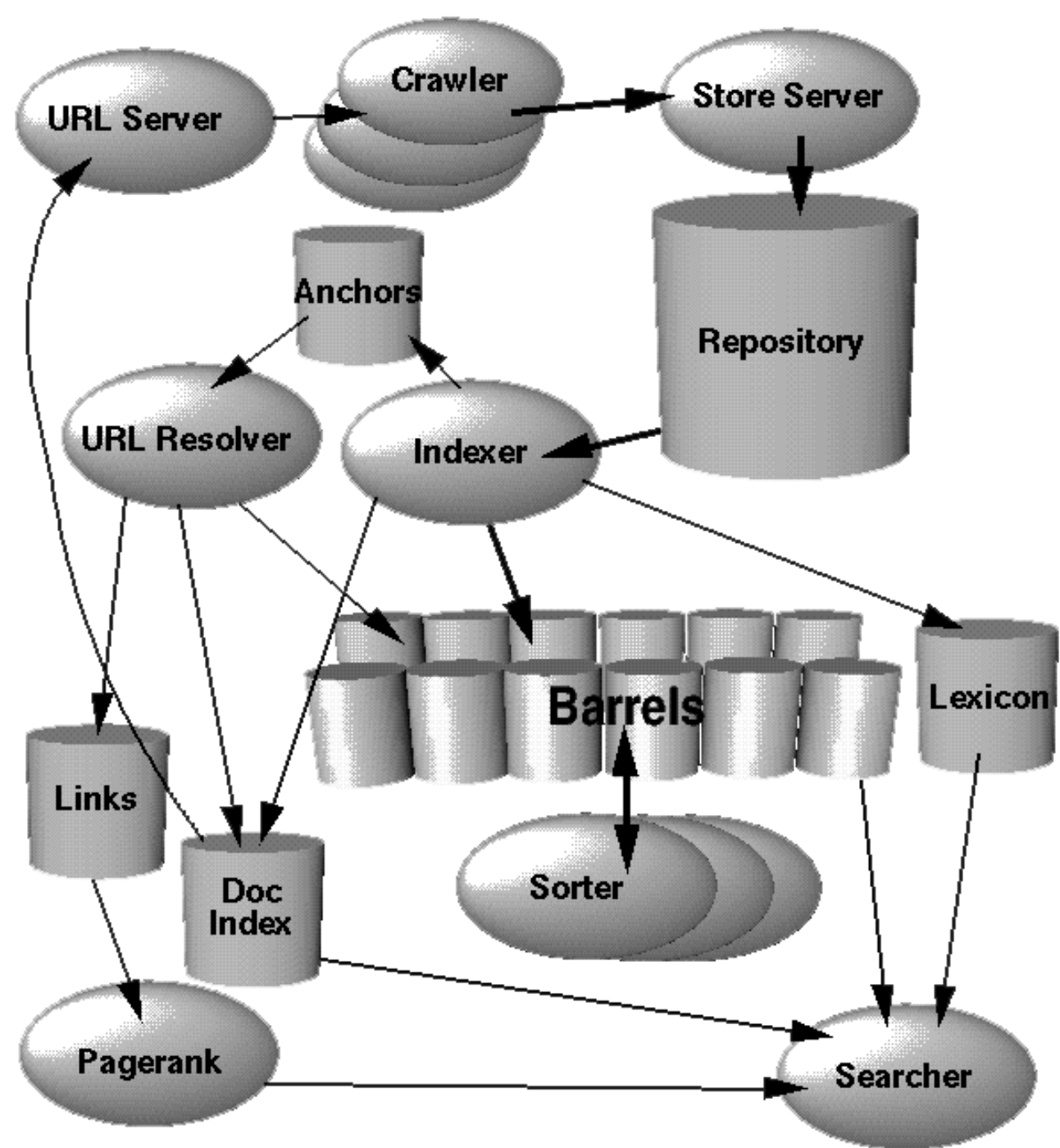
- Apa yang dicache?

➤ Solusi

- Caching banyak level
 - Level Query
 - Level Term

Contoh Lain Arsitektur Search Engine

Search Engine Google



Arsitektur Search Engine

