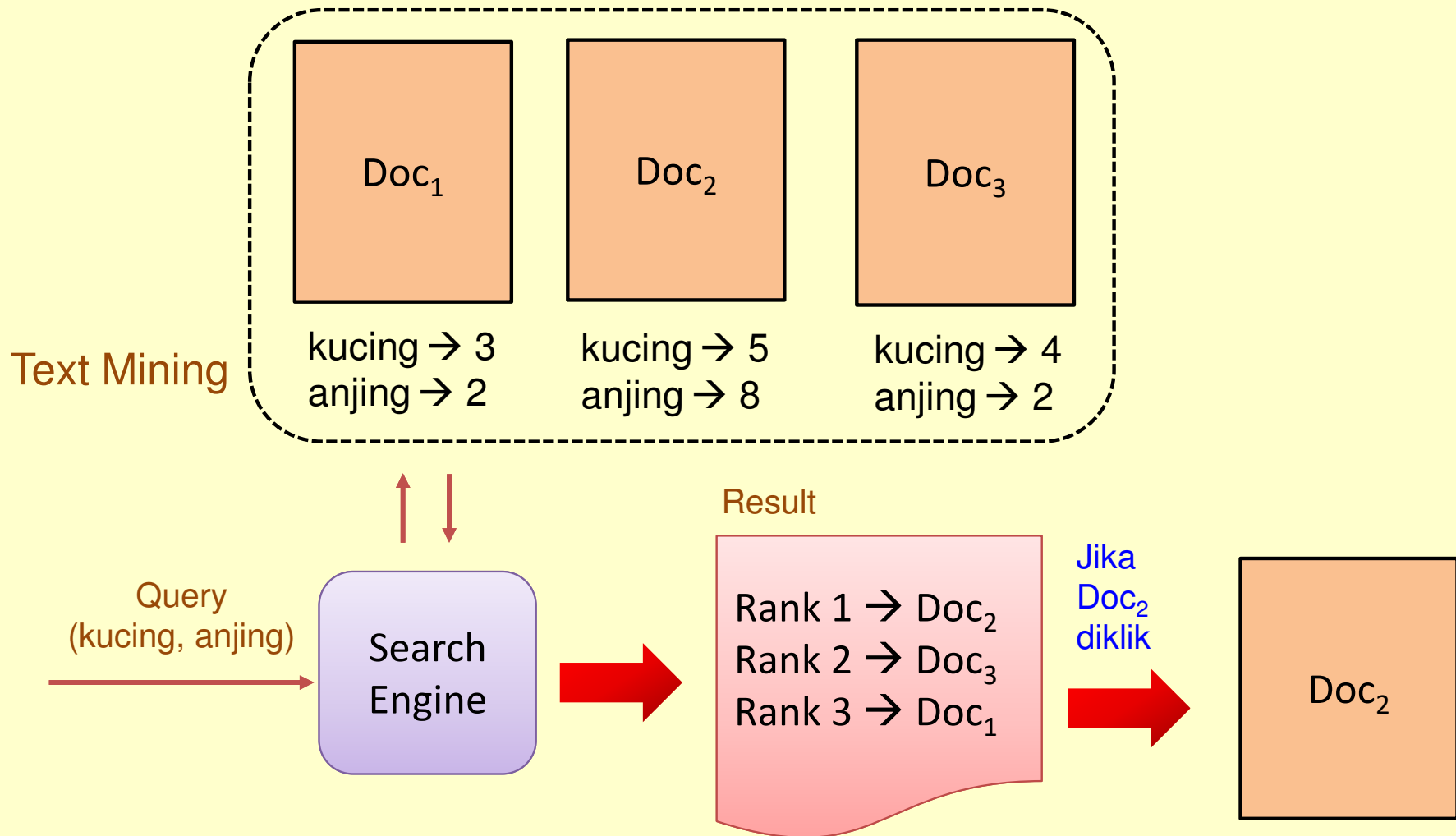# Text Mining & Search Engine

Ali Ridho Barakbah, Entin Martiana, Tri Hadiah
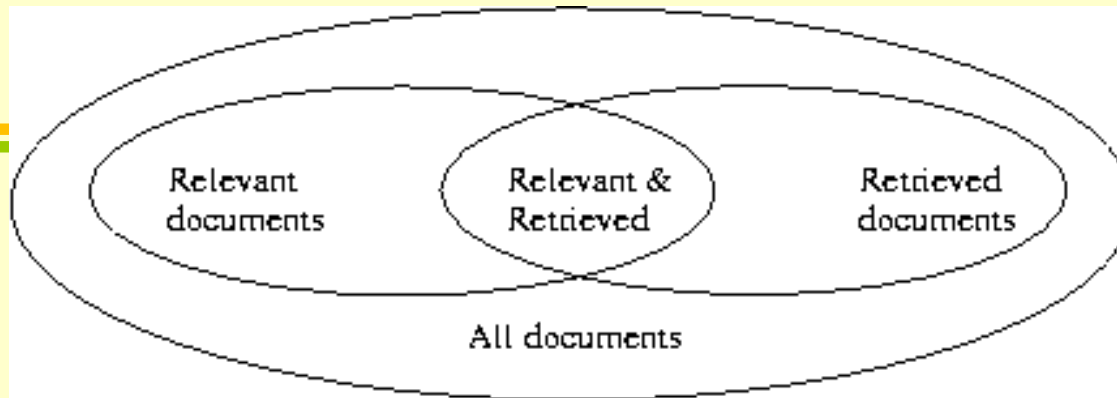
Knowledge Engineering Research Group

Department of Information and Computer Engineering

Politeknik Elektronika Negeri Surabaya

# Text Mining & Search Engine



Text Mining

Doc₁: $kucing \rightarrow 3$, $anjing \rightarrow 2$

Doc₂: $kucing \rightarrow 5$, $anjing \rightarrow 8$

Doc₃: $kucing \rightarrow 4$, $anjing \rightarrow 2$

Query (kucing, anjing) → Search Engine

Result:
Rank 1 → $Doc_2$
Rank 2 → $Doc_3$
Rank 3 → $Doc_1$

Jika $Doc_2$ diklik → $Doc_2$

# Basic Measures for Text Retrieval



- **Precision:** the percentage of retrieved documents that are in fact relevant to the query (i.e., "correct" responses)
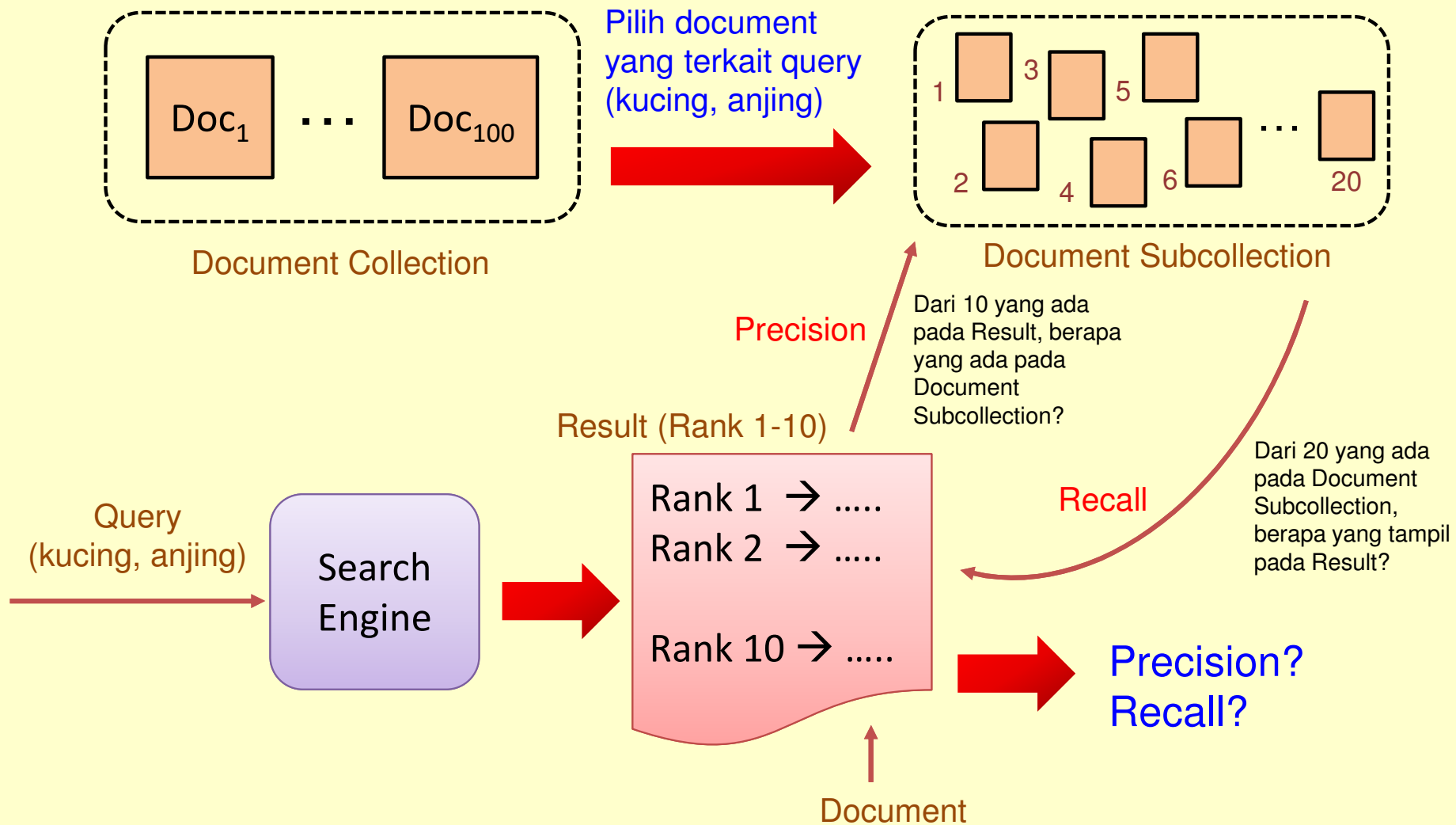
$$precision = \frac{|\{Relevant\} \cap \{Retrieved\}|}{|\{Retrieved\}|}$$

- **Recall:** the percentage of documents that are relevant to the query and were, in fact, retrieved
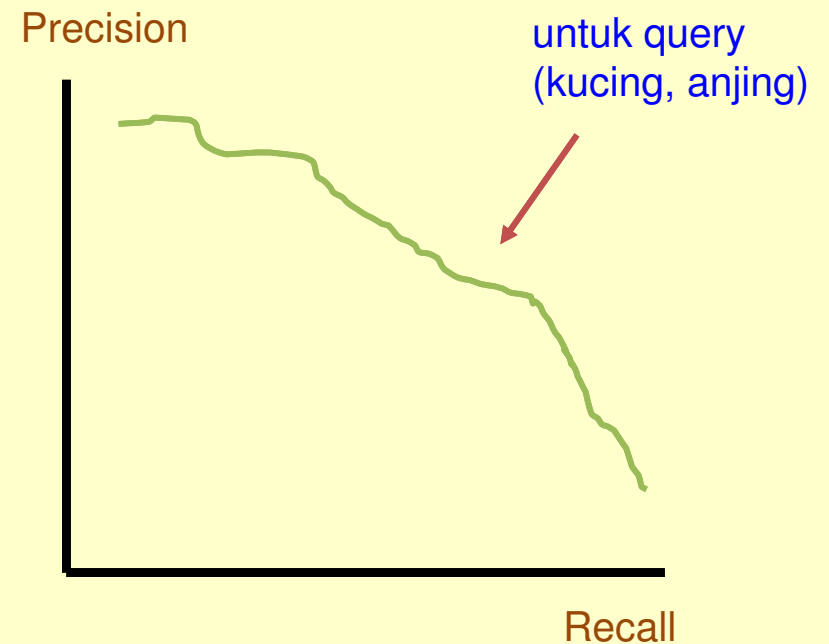
$$recall = \frac{|\{Relevant\} \cap \{Retrieved\}|}{|\{Relevant\}|}$$

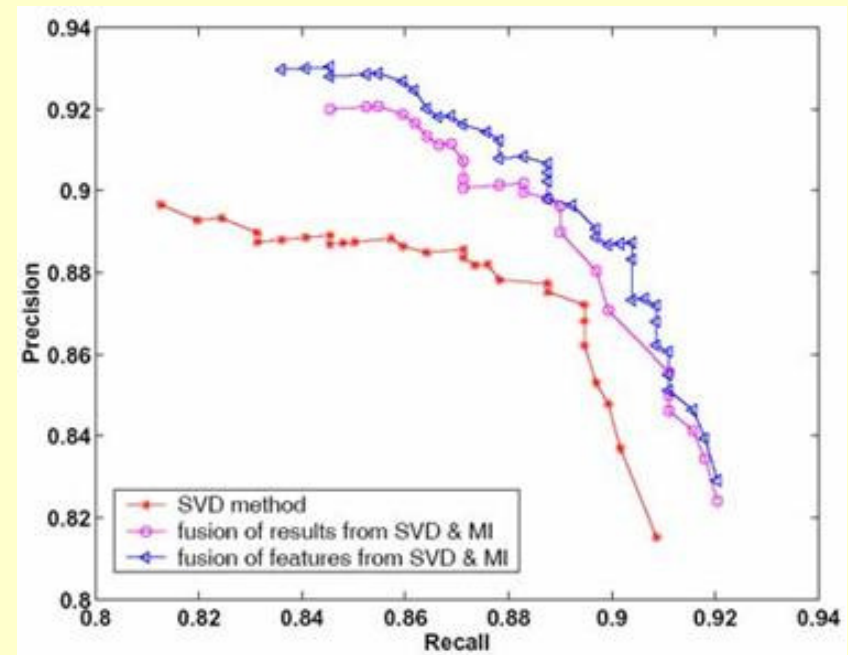Source: Data Mining -Volinsky - 2011 - Columbia University

**Politeknik Elektronika**
**Negeri Surabaya**

# Ilustrasi: Query → kucing anjing

| | | |
|---|---|---|
| Rank 1 | Recall? | Precision? |
| Rank 1-2 | Recall? | Precision? |
| Rank 1-3 | Recall? | Precision? |
| Rank 1-4 | Recall? | Precision? |
| Rank 1-5 | Recall? | Precision? |
| Rank 1-6 | Recall? | Precision? |
| Rank 1-7 | Recall? | Precision? |
| Rank 1-8 | Recall? | Precision? |
| Rank 1-9 | Recall? | Precision? |
| Rank 1-10 | Recall? | Precision? |
| … | Recall? | Precision? |

Precision

untuk query
(kucing, anjing)

Recall

# Precision Recall Curves

# Apa itu Search Engine?

- Software code that is designed to search for information on the World Wide Web. (Wikipedia)

- Programs that search documents for specified keywordsand returns a list of the documents where the keywords were found. (Webopedia)

- Computer software used to search data (as text or a database) for specified information; also : a site on the World Wide Web that uses such software to locate key words in other sites. (Merriam Webster)

# Common Characteristics

- Spider, Indexer, Database, Algorithm

- Menemukan dokumen yang tepat dan menampilkannya sesuai kondisi yang terakhir

- Proses update yang sering terhadap dokumen web pada pencarian dan membuat pemodelan terhadap dokumen

- Berusaha menyajikan hasil yang lebih presisi dibandingkan dengan kompetitor

Source: Saeed El-Darahali, Search Engines & Search Engine Optimization (SEO), 7th World Congress on the Management of e-Business

**Politeknik Elektronika Negeri Surabaya**

Knowledge Engineering (knoWing) Research Group

| Year | Engine | Current status |
|------|--------|----------------|
| 1993 | W3Catalog | Inactive |
|      | Aliweb | Inactive |
| 1994 | WebCrawler | Active, Aggregator |
|      | Go.com | Active, Yahoo Search |
|      | Lycos | Active |
| 1995 | AltaVista | Active, Yahoo Search |
|      | Daum | Active |
|      | Magellan | Inactive |
|      | Excite | Active |
|      | SAPO | Active |
|      | Yahoo! 2008 | Active, Launched as a directory |
| 1996 | Dogpile | Active, Aggregator |
|      | Inktomi | Acquired by Yahoo! |
|      | HotBot | Active (lycos.com) |
|      | Ask Jeeves | Active (rebranded ask.com) |

| Year | Engine | Current status |
|------|--------|----------------|
| 1997 | Northern Light | Inactive |
|      | Yandex | Active |
| 1998 | Goto | Inactive |
|      | Google | Active |
|      | MSN Search | Active as Bing |
|      | empas | Inactive (merged with NATE) |
| 1999 | AlltheWeb | Inactive (URL redirected to Yahoo!) |
|      | GenieKnows | Active, rebranded Yellowee.com |
|      | Naver | Active |
|      | Teoma | Active |
|      | Vivisimo | Inactive |
| 2000 | Baidu | Active |
|      | Exalead | Inactive |
| 2002 | Inktomi | Acquired by Yahoo! |
| 2003 | Info.com | Active |
|      | Scroogle | Inactive |

Source: Wikipedia

| Year | Engine | Current status |
|------|--------|----------------|
| 2004 | Yahoo! Search | Active, Launched own web search (see Yahoo! Directory, 1995) |
|      | A9.com | Inactive |
|      | Sogou | Active |
| 2005 | AOL Search | Active |
|      | Ask.com | Active |
|      | GoodSearch | Active |
|      | SearchMe | Inactive |
| 2006 | wikiseek | Inactive |
|      | Quaero | Active |
|      | Ask.com | Active |
|      | Live Search | Active as Bing, Launched as rebranded MSN Search |
|      | ChaCha | Active |
|      | Guruji.com | Active as BeeMP3.com |
| 2007 | wikiseek | Inactive |
|      | Sproose | Inactive |
|      | Wikia Search | Inactive |
|      | Blackle.com | Active, Google Search |

| Year | Engine | Current status |
|------|--------|----------------|
| 2008 | Powerset | Inactive (redirects to Bing) |
|      | Picollator | Inactive |
|      | Viewzi | Inactive |
|      | Boogami | Inactive |
|      | LeapFish | Inactive |
|      | Forestle | Inactive (redirects to Ecosia) |
|      | DuckDuckGo | Active |
| 2009 | Bing | Active, Launched as rebranded Live Search |
|      | Yebol | Inactive |
|      | Mugurdy | Inactive due to a lack of funding |
|      | Goby | Active |
|      | NATE | Active |
| 2010 | Blekko | Active |
|      | Cuil | Inactive |
|      | Yandex | Active, Launched global (English) search |
|      | Yummly | Active |
| 2011 | Interred | Active as Interredu |
|      | Yandex | Active, Launched Turkey search |
| 2012 | Volunia | Active |
|      | Interredu | Active |
|      | Open Drive | Active, cloud file search |
| 2013 | iStella | Active |
|      | Aoohe | Active |

Source: Wikipedia

Politeknik Elektronika
Negeri Surabaya

Knowledge Engineering
(knoWing) Research Group
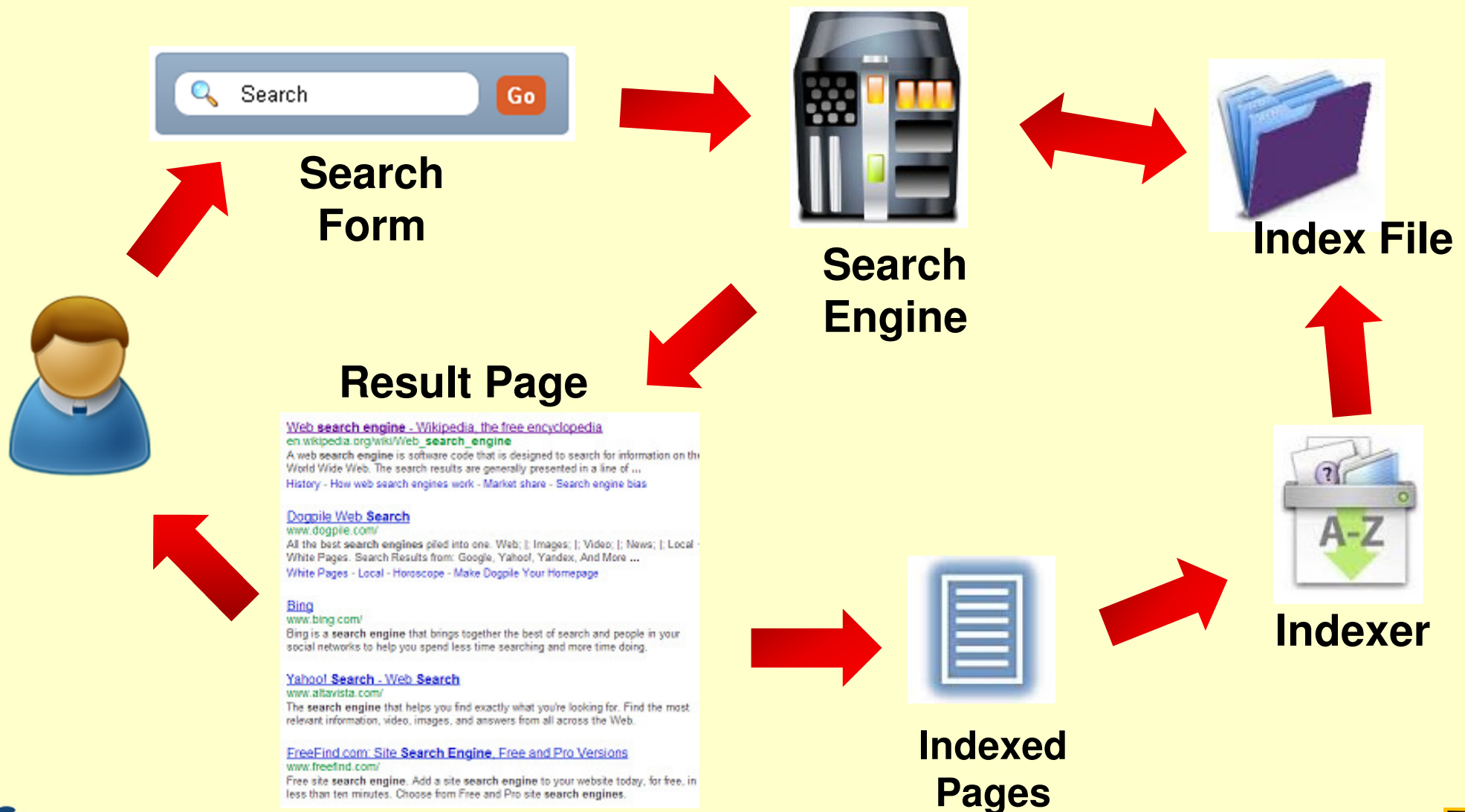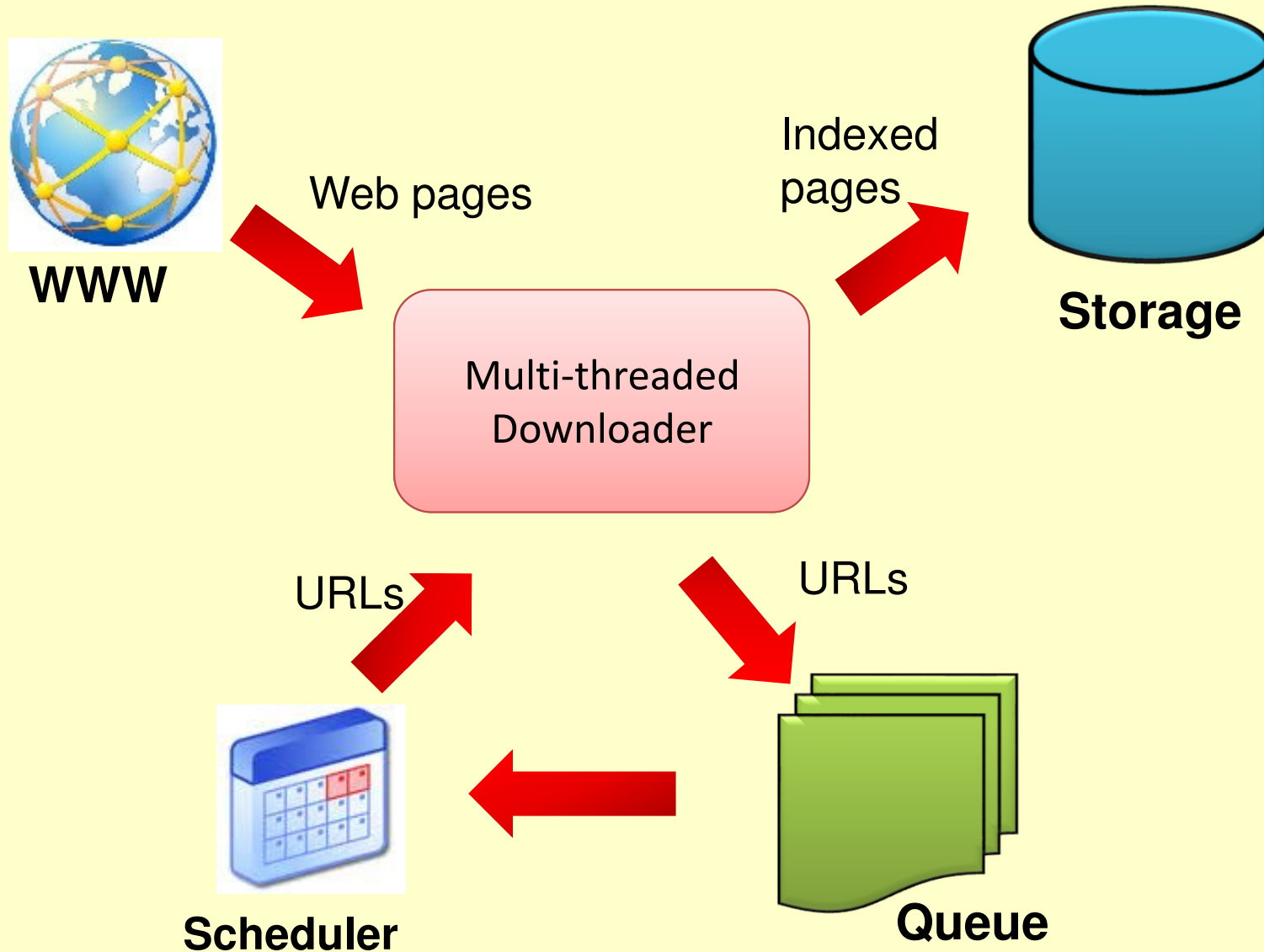
# Bagaimana Search Engine Bekerja?

- Spider melakukan crawling halaman-halaman web untuk menemukan dokumen-dokumen baru, biasanya dengan mengikuti hyperlinks dari web yang sudah ada di database

- Search engine melakukan indexing terhadap halaman web dan menambahkannya ke dalam database. Ada proses update secara berkala.

- Search engine melakukan pencarian pada database berdasarkan query yang dimasukkan oleh user (bukan langsung pencarian pada halaman web)

- Search engine melakukan ranking dari hasil pencarian dokumen dengan menggunakan algoritma tertentu
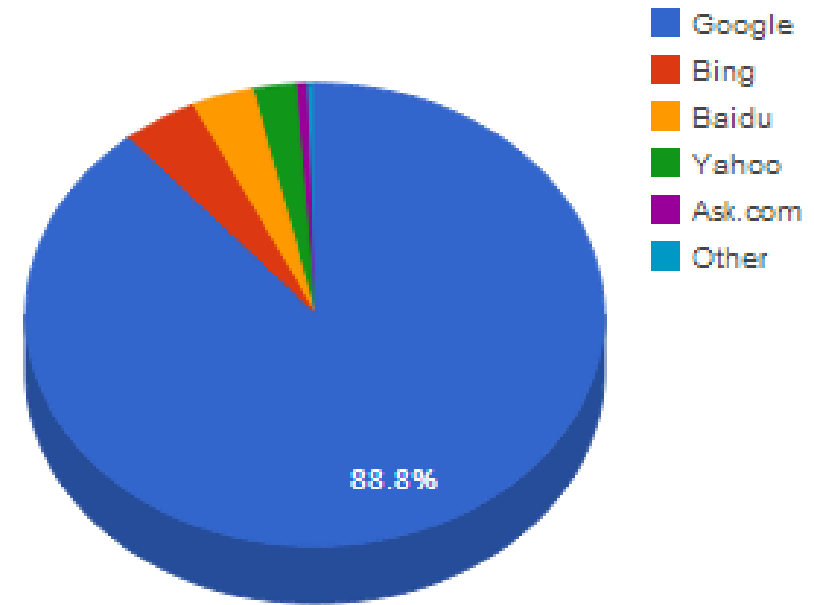
# Bagaimana Web Crawler bekerja?



WWW

Web pages

Multi-threaded Downloader

Indexed pages

Storage

URLs

Queue

URLs

Scheduler

# Market Share dari Search Engine

Source:
www.karmasnack.com/about/search-engine-market-share/

**Global:**



| Global: | |
|---------|---------|
| Google | 88.8% |
| Bing | 4.2% |
| Baidu | 3.5% |
| Yahoo | 2.4% |
| Ask.com | 0.6% |
| Other | 0.5% |