

# LC029 정보검색

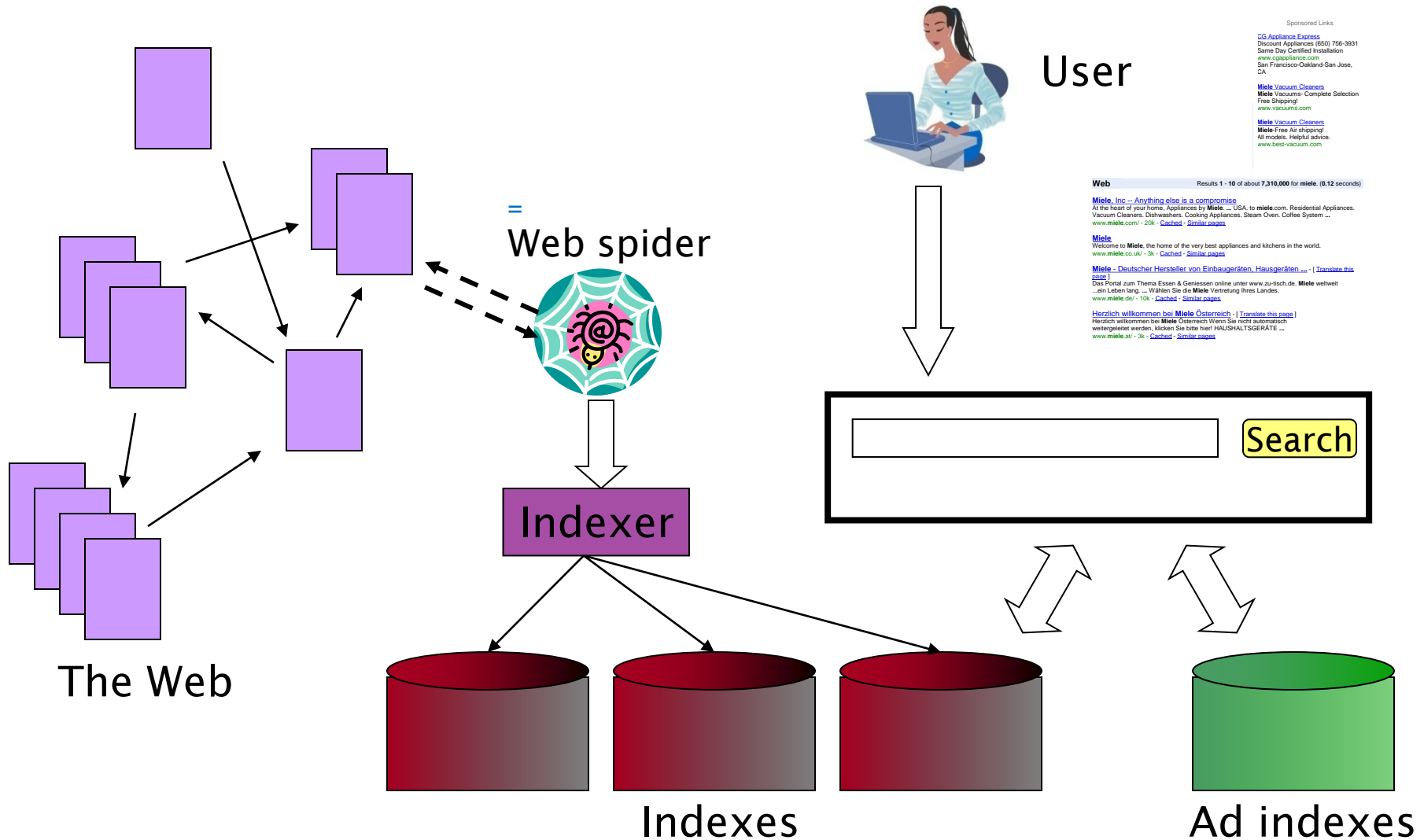
2022 11 17

## Chapter 19 : Web Search Basics

Web Search Characteristics  
Advertising as Economic Model & Spam  
Index Size Estimation  
Duplicate Documents

# Web Search Characteristics

# Web Search Basics



# Characteristics of Web document collection

---

## ■ Creation

- Decentralized : no design and no co-ordination.
- Distributed and democratized content creation.
- Diversity of backgrounds and motives of its participants.

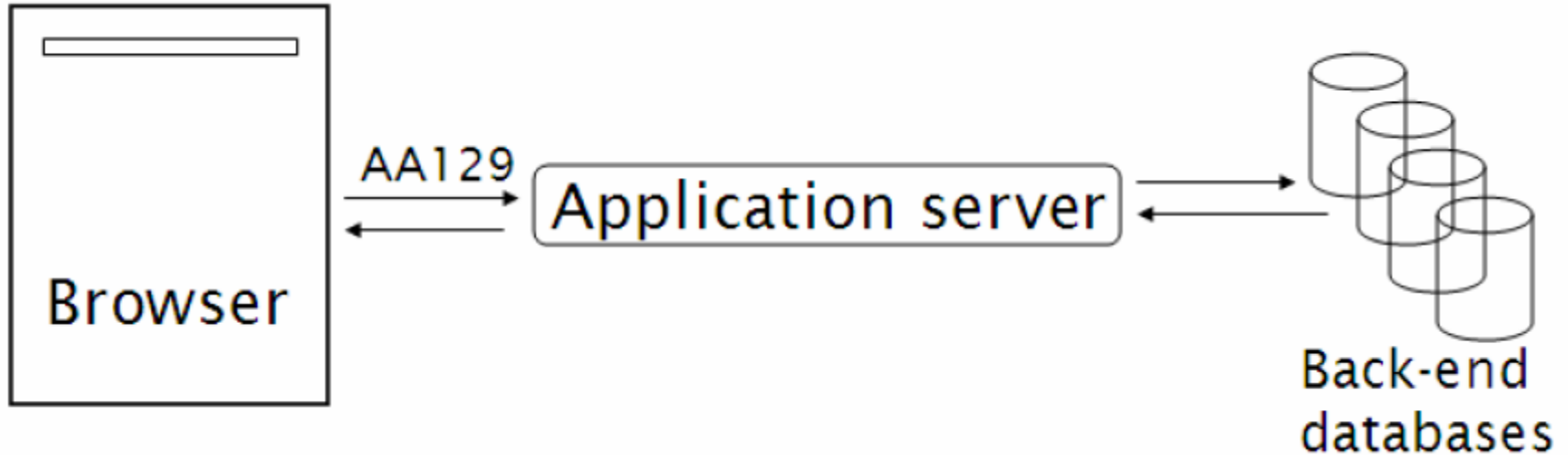
## ■ Contents

- Includes truth, lies, obsolete information, contradictions.
- In many languages (many variations in grammar and style)  
-> require different stemmers and linguistic operations
- Can be either static or *dynamically generated*.

static

# Characteristics of Web document collection

---



dynamic generation of web pages

: 가

가

# Characteristics of Web document collection

---

## ■ Structures

- Unstructured (text, html, ...)
- Semi-structured (XML, annotated photos)
- Structured (Databases)

## ■ Scale

- Much larger than previous text collections.
- By 1995, Altavista has crawled and indexed 30 M pages.
- In 1995, the volume of web pages was doubling every few months.
- The growth is now slowed down, but still expanding.

# Web Information Retrieval

---

- Web Information Retrieval
  - Web document collections are useful when they are discoverable. 가 .
  - **Keyword-based search** : Altavista, Excite, Infoseek
    - Based on inverted index and ranking mechanism ,
  - **Taxonomy-based search** : Yahoo 가 가
    - Convenient and intuitive for finding web pages.
    - For accurate and consistent classification, manual classification is needed.
    - Manual classification entails significant human efforts.
    - Mismatch between users and editors of the taxonomy. 가
    - As the size of taxonomy grows, it is hard for users to find information.



# Web Information Retrieval

---

- First generation of Web Information Retrieval
  - Used classical search techniques we have learned.
  - For indexing, query processing and ranking **at web scale**, they need to harness tens of machines together.
  - The **quality and relevance** of web search results are not satisfactory due to the idiosyncrasies of content creation on the web.
  - **New ranking** and **spam-fighting** techniques have invented to improve the quality of search results.
  - Need to measure the **authoritativeness** of a document based on cues such as which website hosts it.

가

가

# Web Search User Experience

# Web Search User Experience

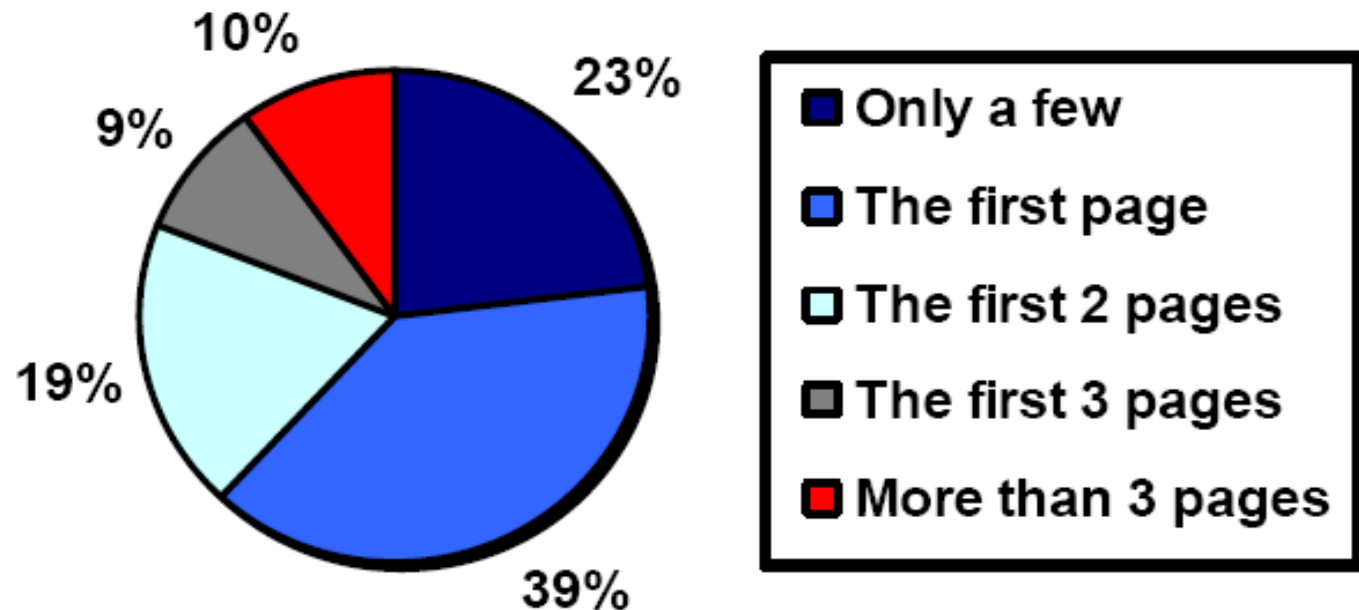
---

- Traditional IR : 가가 ,
  - Users were typically professionals.
  - They are trained in the art of phrasing queries.
- Web Search : 가가 ,
  - Users do not know about the heterogeneity of web content.
  - They do not know the syntax of query language, the art of phrasing queries.
  - Average number of keywords is 2-3.
  - They seldom use operators such as AND, OR and wildcards.

# How far do people look for results?

---

When you search something, how many entries do you typically review before clicking one?

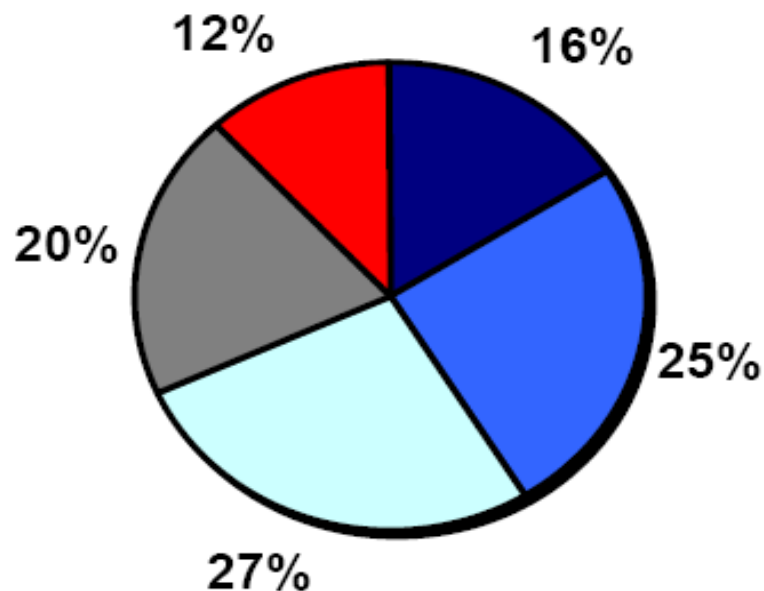


(Source: [iprospect.com](http://iprospect.com) WhitePaper\_2006\_SearchEngineUserBehavior.pdf)

# How far do people look for results?

---

If you don't find what you are looking for, when do you revise your search or move on to another search engine?



- After reviewing the first few entries
- After reviewing the first page
- After reviewing the first 2 pages
- After reviewing the first 3 pages
- After reviewing more than 3 pages

# User Query Needs

---

- Web search queries are categorized into three groups.
  - Informational Low hemoglobin
    - Want **to learn** about something (40% ~ 65%).
    - **Try to assimilate information from multiple web pages.**
  - Navigational 가 Sungshin Women's University
    - Want **to go** to that page (15% ~ 25%).
    - **Measure of user satisfaction is precision at 1.**
  - Transactional ,
    - Want **to do** something (20% ~ 35%).
      - Purchasing a product Apple ipod
      - Downloading a file Mars surface images
      - Accessing a service Road condition of Seoul
    - **Search engine should return results listing those services.**

# User Query Needs

---

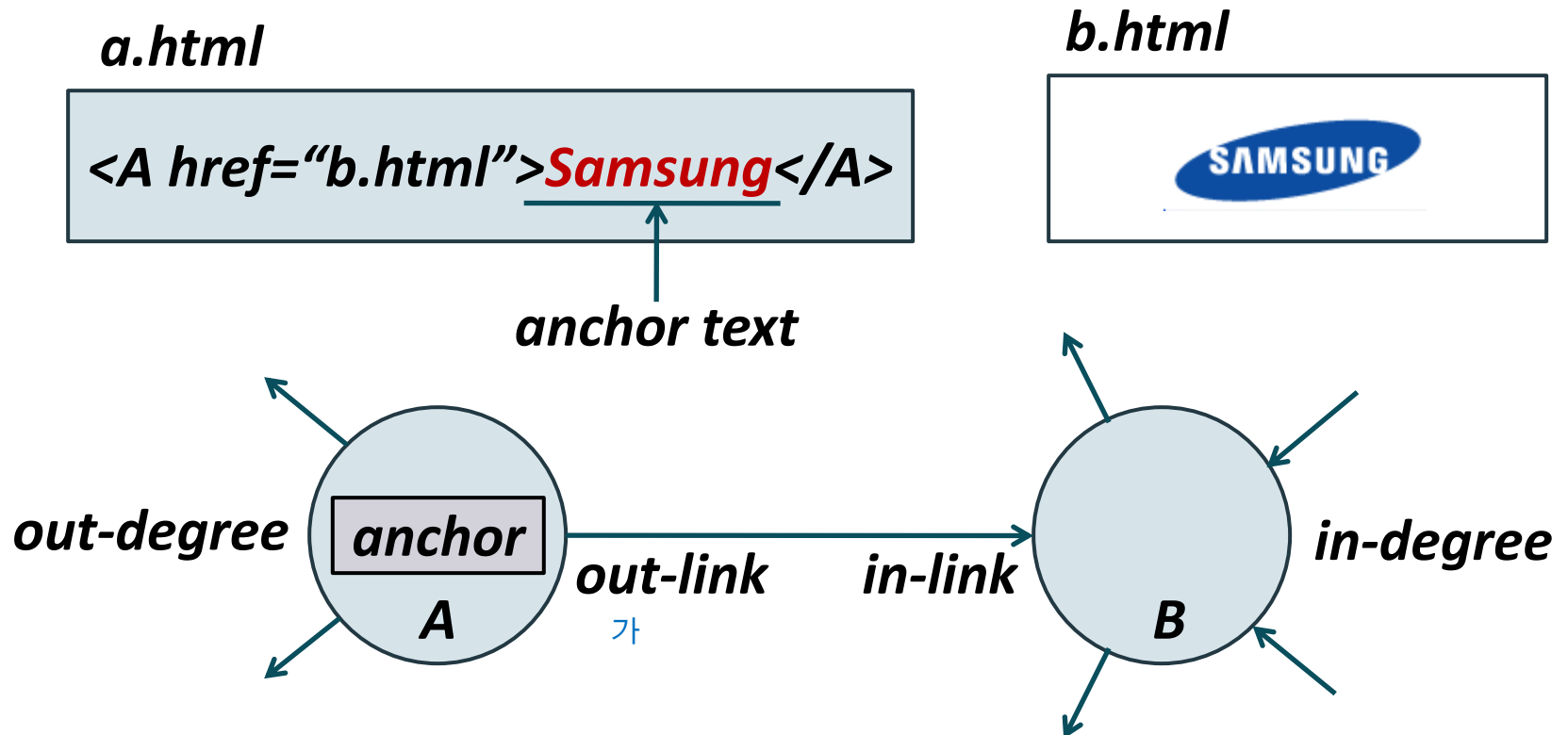
- Other than that, users want to: :
  - Find a good hub. 가 Car rental Brasil
  - Do an exploratory search to see what is going on.

# Web Graph



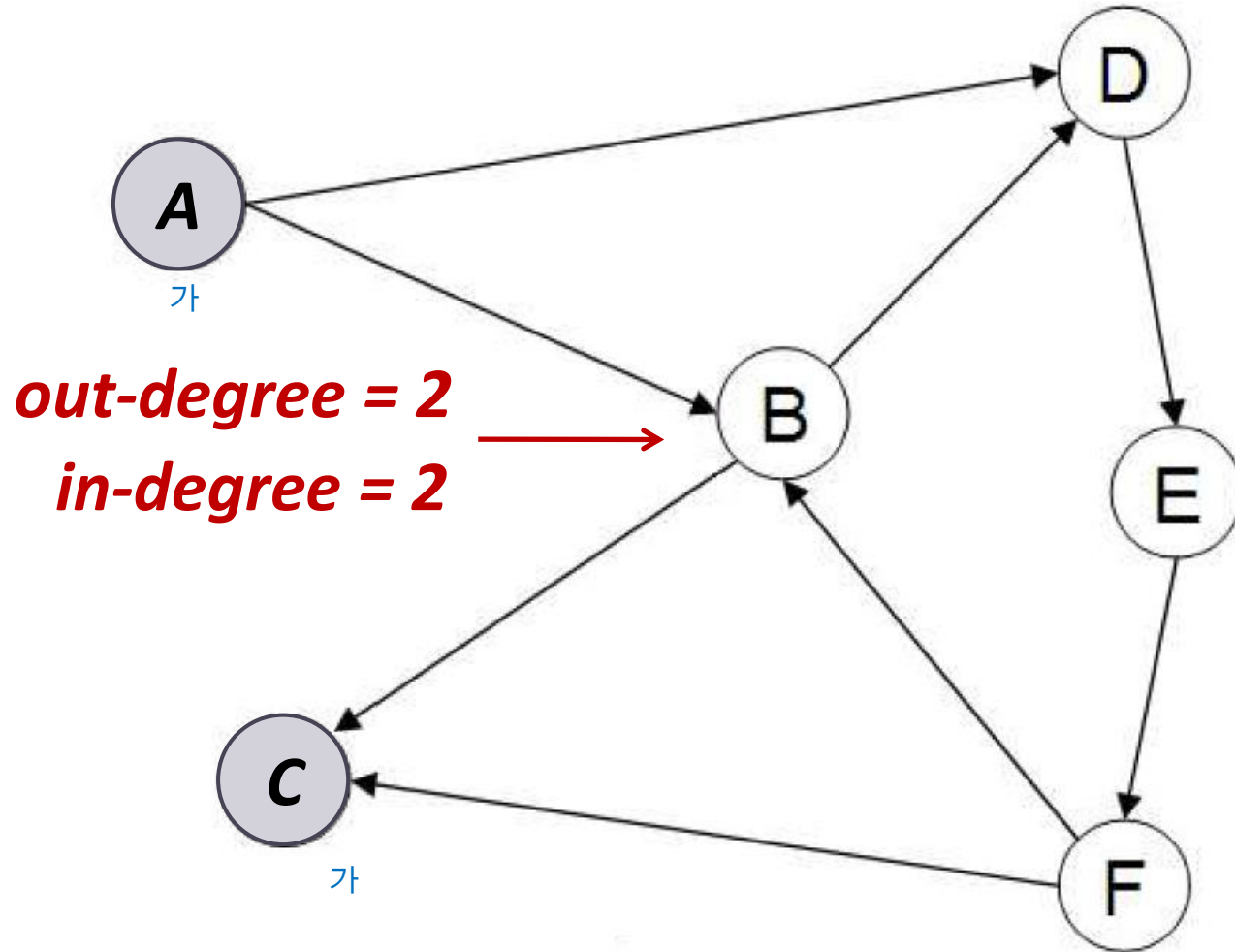
# Web Graph

- Web documents have hyperlinks between them as a directed graph.



# Web Graph

---



# Web Graph

- The links are not randomly distributed.
  - The distribution of the number of links *into* a web page does not follow the Poisson distribution.
  - The distribution follows **power law**:  
The number of web pages with in-degree  $i \propto \frac{1}{i^\alpha}$  where  $\alpha$  is typically 2.1.

**when  $i = 1$ ,  $1/1^{2.1} = 1.00$**

**when  $i = 2$ ,  $1/2^{2.1} = 0.23$**

**when  $i = 3$ ,  $1/3^{2.1} = 0.10$**

**when  $i = 4$ ,  $1/4^{2.1} = 0.05$**

**when  $i = 5$ ,  $1/5^{2.1} = 0.03$**

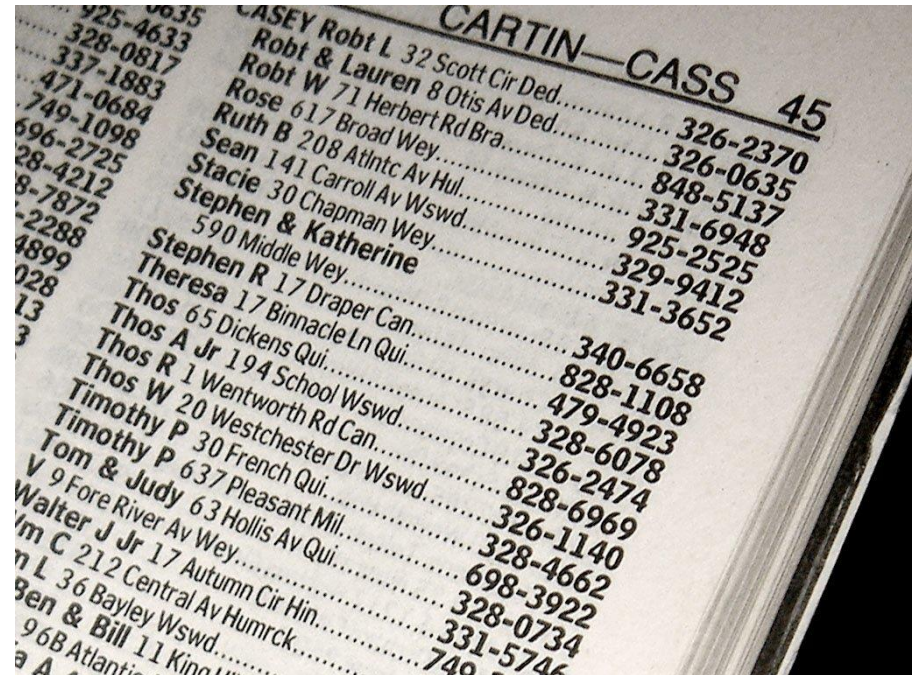
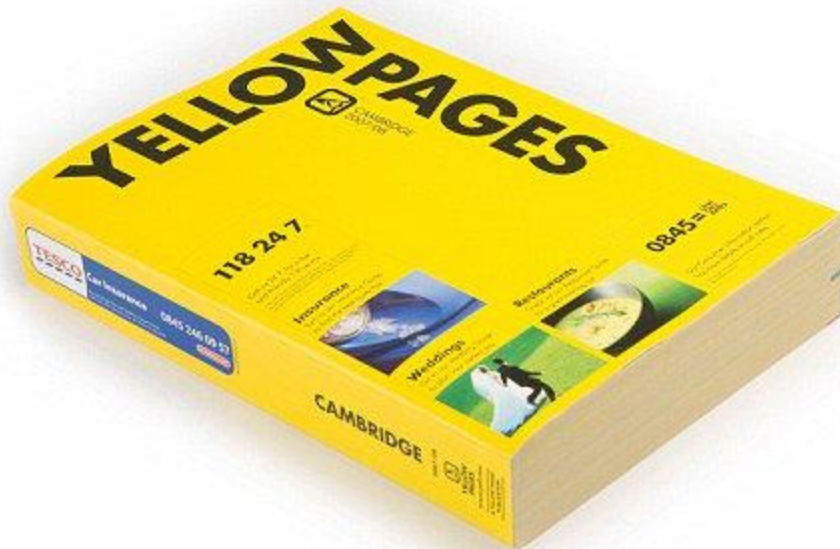
가

가

# Advertising as Economic Model

# Paid Inclusion Model

- The Yellow Pages
  - Companies pay for larger / darker fonts.
  - It is fair because they pay for that.



# Paid Inclusion Model

- The Yellow Pages
  - Companies pay for larger / darker fonts.
  - It is fair because they pay for that.



# Paid Inclusion Model

---

- The Yellow Pages

- Some companies use some tricky method to list their name early in the list. ( )
  - **aaaaaaCheapestPhoneCompanyaaaaa**
  - They do not pay for that.
- Web search engine followed the Yellow Pages' paid inclusion model.

CheapestPhoneCompany

: 가 가 "가 " ( )

# Advertisement

---

- Early Stage
  - **Banner advertisement** at popular websites (MSN, CNN, ...)
  - The purpose of the advertisement was to convey a positive feeling about their **brand**.
  - The advertisement was priced on CPM (cost per mille).
- Later
  - As the goal of the advertisement was to induce a **transaction**, the pricing model was changed to CPC (cost per click).



# Sponsored Search : GOTO

---

- Goto was not a search engine in the traditional sense.
  - Your search ranking depended on how much you paid.
  - When the user clicked on a result page, Goto was compensated by the corresponding advertiser.
- Goto users are actively expressing an interest and intent related to the query term.
  - A user typing **iphone** is more likely to buy it than simply browsing news about **iphone**.
  - Auction for keywords: **casino** was expensive!

:

# Sponsored Search : GOTO

---

- Goto was morphed into Overture, and finally acquired by Yahoo!
  - Meanwhile Goto/Overture's annual revenues were nearing \$1 billion.
- Google added paid-placement ads to the side, independent of search results.
  - Pure search engine + **Sponsored search** engine

# Sponsored Links

Google   [Advanced Search](#)

Web [+ Show options...](#) Results 1 - 10 of about 45,500,000 for **samsung tv**. (0.28 seconds)

[Television | SAMSUNG](#)  
Plasma TV. Type main TV. Which **Samsung TV** is Right for You? **Samsung** TVs offer world-class picture quality, design and energy efficiency. Find the TV that's ...  
[www.samsung.com/uk/consumer/tv-audio-video/.../index.idx?... - Cached](http://www.samsung.com/uk/consumer/tv-audio-video/.../index.idx?...)

[Plasma TV - Televisions | SAMSUNG](#)  
Get more out of your **Samsung TV** with a host of connectivity options. ...  
[www.samsung.com/us/consumer/tv-video/televisions/...tv/index.idx?... - Cached](http://www.samsung.com/us/consumer/tv-video/televisions/...tv/index.idx?...)

[TV | SAMSUNG](#)  
LED; Plasma; LCD. **Samsung** LCD and Plasma Series; **Samsung** ECO TV; LED TV ...  
[www.samsung.com/ca/consumer/type/type.do?group=tv&type=tv](http://www.samsung.com/ca/consumer/type/type.do?group=tv&type=tv)

[More results from samsung.com »](#)

[Samsung Lcd Tv - Televisions - Compare Prices, Reviews and Buy at ...](#)  
**Samsung** Lcd Tv - 170 results like the **Samsung** LN52A650 52 in. HDTV LCD TV, **Samsung** LN-46A650A 46" LCD HDTV, **Samsung** LN52B750 52" LCD HDTV, **Samsung** UN55B8500 ...  
[www.nextag.com/samsung-lcd-tv/stores.html - Cached - Similar](http://www.nextag.com/samsung-lcd-tv/stores.html)

[Samsung TV Reviews, Samsung Television Reviews & Ratings](#)  
We've analyzed price, features, and reviews of **Samsung** TVs to find the best values. Also find quick links to the most useful user reviews for all **Samsung** TV ...  
[www.retrevo.com/samples/Samsung-TV.html - Cached - Similar](http://www.retrevo.com/samples/Samsung-TV.html)

**Ads** →

**Sponsored Links**

[Sale on White 32" LCD TV](#)  
Take 30% Off Today Only!  
Brand New. Free Shipping  
[www.eBay.com/32\\_tv\\_white+sale](http://www.eBay.com/32_tv_white+sale)

[Top 10 LCD Televisions](#)  
Compare Prices & Save - Shop Smart!  
PC World Shopping: Tv Lcd  
[LCD-TVs.PCWorld.com](http://LCD-TVs.PCWorld.com)

[Samsung Tv](#)  
Find a Solution for any  
Support problem Easily  
[www.Fixya.com](http://www.Fixya.com)

[See your ad here »](#)

← **algorithmic search results**

# Sponsored Links

Google

samsung tv

All Images News Shopping Videos More Settings Tools

About 1,500,000,000 results (0.61 seconds)

**Samsung All TVs - Explore 8k, 4k & UHD Smart TVs | Samsung US**  
<https://www.samsung.com> > Home > Televisions Home Theater > Tvs > All Tvs ▼  
Compare All types of TV models by **Samsung**. QLED, UHD, Full HD TVs are available in various sizes and equipped with smart features and big screen.

Top stories

**Samsung vs LG TV: which TV brand is better?**  
TechRadar  
1 day ago

**Walmart drops great deals on Samsung 4K smart TVs for Memorial ...**  
Digital Trends  
7 hours ago

**Best 4K TV deals this week: Save on Samsung, LG, and more**  
Mashable  
15 hours ago

**See samsung tv** Sponsored ⓘ

**삼성 55인치 UHD UN55NU6900 스마트 TV**  
₩819,000  
11번가  
Free shipping

**삼성55인치LED-TV UN55K5110BF 인터넷 +IPTV가입**  
₩260,000  
위메프  
Free shipping

→ More on Google Google is not a party to the product sale

*Ads →*

*algorithmic search results ←*

→ More for samsung tv

# Sponsored Links

**NATE**   검색 통합검색 S 시맨틱 이미지 동영상 사람검색

---

바로가기 [애플 마이팟](http://www.apple.com/kr/...) <http://www.apple.com/kr/...>

---

**스폰서링크** AD

- [애플iPod공인대리점 맥이샵](http://www.maceshop.co.kr) - 애플 iPod과 컴퓨터 및 액세서리 파격행사최대80% 무이자할부, 매장구입도환영.  
<http://www.maceshop.co.kr>
- [ipod 전문쇼핑몰 키스맥몰](http://mall.kissmac.com) - ipod액세서리 70% 할인, 파우치 공짜, 독세이버 500원, 케이블러 1천원,  
<http://mall.kissmac.com>
- [재밋는 쇼핑2.0 오픈베이](http://www.openbay.co.kr) - 마이팟, 클래식120G 26500원, 셔플5 4250원, 터치8G 3650원,  
<http://www.openbay.co.kr>

Spam

# Trouble with Sponsored Links

---

- To rank highly, it costs money.
- An alternative is ***Search Engine Optimization (SEO)***:
  - Tuning your web page to rank highly in the **algorithmic search results** for selected keywords.
  - It is an alternative to paying for placement, thus, intrinsically a marketing function.
  - Performed by companies (**Search Engine Optimizers**) for their clients.
  - Some perfectly legitimate, some very shady.

# Search Engine Optimization : **Keyword Stuffing**

---

- First generation Search Engines relied heavily on  $tf \cdot idf$ . Web pages with high  $tf$  would rank highly.
  - This led to the first generation of spam.
  - Sophisticated spammers rendered the repeated terms in the same color as the background.
    - Repeated terms got indexed by indexers.
    - But not visible to humans on browsers.
  - Thus, pure word density cannot be trusted as an IR signal.

$tf \cdot idf$

$\cdot idf$

$tf$

$\cdot$



# Search Engine Optimization : **Keyword Stuffing**

---

- They also use meta-tags.

```
<meta name="description" content="... London  
hotels, hotel, holiday inn, hilton, discount, booking,  
reservation, mp3, britney spears, ..." />
```

- **Meta tags** are snippets of text that describe a page's content.
- Meta tags don't appear on the page itself, but only in the page's source code.
- Meta tags are essentially little content descriptors that tell Search Engines what a web page is about.

# As a result, ...

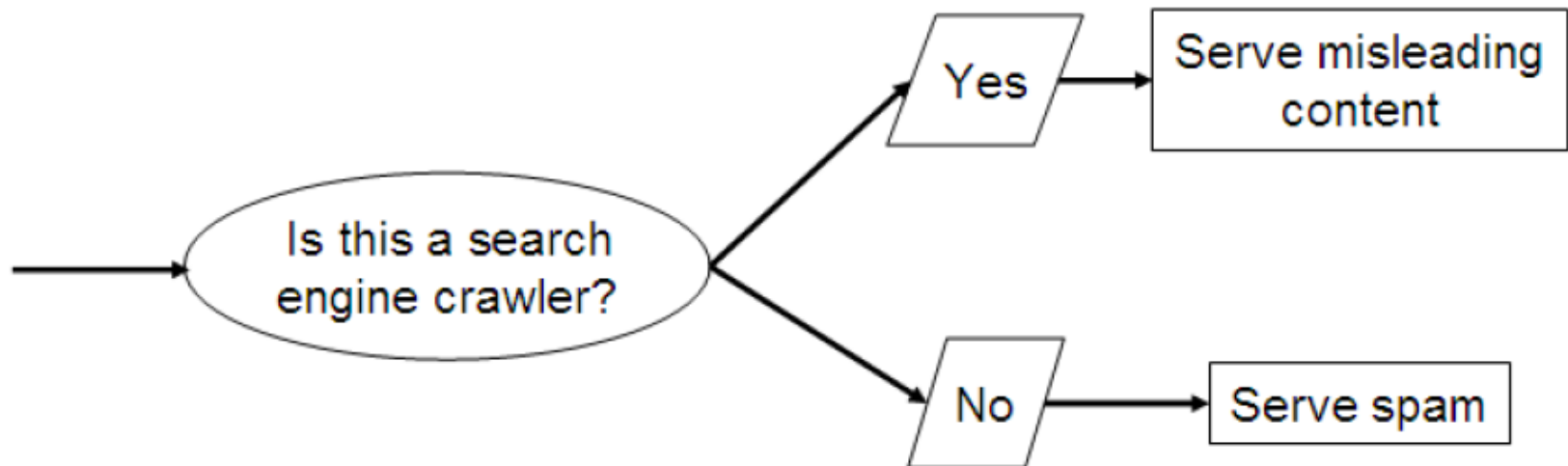
---

- Search Engine became sophisticated enough to screen out a large number of repetitions of particular keywords.
- Spammers responded with a richer set of spam techniques.
  - Cloaking
  - Doorway Page
  - Click Spam

# Search Engine Optimization : Cloaking

---

- Serve fake content to search engine spider.
  - Indexed by the search engine under misleading keywords.
- When users search for these keywords, they receive different web pages.



# Search Engine Optimization : Doorway Page

---

- A doorway page contains text and meta data carefully chosen to rank highly on selected keywords.
- When a browser request the doorway page, it is redirected to a more commercial page.

```
<HTML>
<HEAD>
<meta http-equiv="refresh" content="1" ; url="fake.com" />
</HEAD>
<BODY>
... text and metadata carefully chosen to rank highly on
    selected keywords ...
</BODY>
</HTML>
```

# Search Engine Optimization : Click Spam

---

- No universally accepted definition of click spam.
- Example
  - Repeatedly clicking on *sponsored search results* to exhaust the advertising budget of a competitor.

⋮

# Index Size Estimation

# Index Size Estimation

---

- What is the size of web?

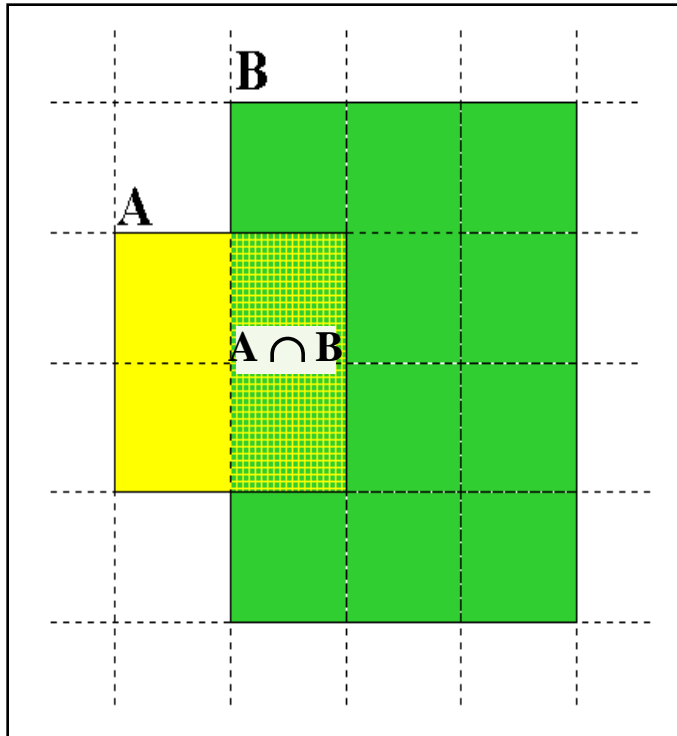
- The web is really infinite.
- Static web  
contains duplication, mostly due to mirroring (~40%).
- Dynamic web  
e.g. calendar, HTTP 404

HTTP 404 - 파일을 찾을 수 없음  
Internet Explorer

- What can we measure?

- The static web pages are whatever search engines have to index.
- To know the coverage of a search engine relative to another, we measure the **relative index sizes**.

# Relative Index Sizes of Two Engines



**Sample** choose URLs randomly from A.  
**Check** if contained in B.

$$\rightarrow A \cap B = (1/2) \times \text{Size A}$$

**Sample** choose URLs randomly from B.  
**Check** if contained in A.

$$\rightarrow A \cap B = (1/6) \times \text{Size B}$$

$$(1/2) \times \text{Size A} = (1/6) \times \text{Size B}$$

$$\therefore \text{Size A} / \text{Size B} = (1/6) / (1/2) = 1/3$$

$$\text{Size A} = 1/3 \times \text{Size B}$$

$$\text{Size B} = 3 \times \text{Size A}$$

**Each test involves:** (i) Sampling URLs (ii) Checking



# Relative Index Sizes of Two Engines

---

- Ideal Strategy

- Generate a random, *uniformly-distributed* URL.
- Check for containment in the index of each engine.

가

- Problem

- Random, uniformly-distributed URLs are hard to find!
- It is enough to generate a *random URL contained in a given Engine*.

가

- Approach 1: Random Searches
- Approach 2: Random IP Addresses
- Approach 3: Random Queries

4 : IP IP 4

# Sampling URLs

---

- Approach 1: Random Searches
  - Send a random search from a search log.
    - Search log is an accumulation of all search queries of a work group.
    - This may include the *bias* from the types of searches made by the work group.
  - Pick a random page from the search results.
  - Check if the page is contained in the other engine.

# Sampling URLs

---

- Approach 2: Random IP Addresses
  - Generate a random IP address. [IP](#)
  - Send a request to a *web server of that address*.
  - Collect all pages at that server.
  - Check if those pages are contained in each engine.

# Sampling URLs

---

- Approach 3: Random Queries
  - Pick a set of random terms from dictionary.
    - **Lexicon: 400,000+ words from a web crawl**
  - Form a random query with two or more terms connected conjunctively. ( , )
  - Do search with the random conjunctive query.
  - Pick a page  $p$  at random from the top 100 returned results.
  - Check if the page  $p$  is contained in the other engine.

: 가

# Duplicate Documents

# Duplicate Documents

---

- The web is full of duplicated content.
  - Mirroring for reliability
  - 40% of the web pages are duplicates of other pages.
  - No need to index multiple copies of the same pages.
    - Save storage space
    - Reduce processing overheads

가

# Duplicate Documents

---

## ■ Duplication

- Can be detected with a ***fingerprint***, a succinct (say 64-bit) digest of the characters on a page.
- If two fingerprints are equal, we test whether the pages are **really equal**.

## ■ Near-Duplication

- There are many cases of near duplication on the web.
- Compute similarity of two documents.
  - **Jaccard Coefficient** can be used.
- Use similarity threshold to detect near duplicates.
  - e.g. Similarity > 80% => Documents are “near duplicates”

# Jaccard Coefficient

---

- Measurement of the overlap of two sets  $A$  and  $B$ .
  - $\text{Jaccard}(A, B) = |A \cap B| / |A \cup B|$
- Always assigns a number between  $0$  and  $1$ .
  - $\text{Jaccard}(A, A) = 1$
  - $\text{Jaccard}(A, B) = 0$  if  $A \cap B = \emptyset$
- $A$  and  $B$  don't have to be the same size.
- Used to measure the similarity of  $A$  and  $B$ .



# Jaccard Coefficient: Scoring Example

---

- For each of the two documents below, calculate Jaccard Coefficient of them.
  - Document 1: caesar died in march
  - Document 2: the long march in Korea

$$\text{Jaccard}(d1, d2) = 2/(4 + 5 - 2) = 1/7 = 0.143$$

# Computing Similarity

---

- Are the documents duplicates of each other?

***d1***

***can a can can a can? a  
rose is a rose. she sells  
sea shells.***

***d2***

***she sells a rose. a rose  
is a can. can a can can  
sea shells.***

- $\text{Jaccard}(d1, d2) = 15 / (15 + 15 - 15) = 1.0$   
Therefore, d1 is a duplicate of d2, or vice versa.
- In reality, they are not!!

:

# Computing Similarity

- Use  $k$ -gram of terms.

$k$   
4

***d1***

***can a can can a can? a  
rose is a rose. she sells  
sea shells.***

***can a can can  
a can can a  
can can a can  
can a can a  
a can a rose  
can a rose is***

***...***

***d2***

***she sells a rose. a rose  
is a can. can a can can  
sea shells.***

***she sells a rose  
sells a rose a  
a rose a rose  
rose a rose is  
a rose is a  
rose is a can***

***...***

# Computing Similarity

---

- Shingling of a document
  - $k$ -shingles of a document  $d$  are defined to be a set of all  $k$ -grams in  $d$ . k
  - ***a rose is a rose is a rose*** (assume  $k = 4$ )  
a\_rose\_is\_a  
rose\_is\_a\_rose  
is\_a\_rose\_is  
a\_rose\_is\_a
- Two documents are near duplicates if the sets of shingles are nearly the same.

# Computing Similarity

---

- $S(d_j)$  : the set of shingles of document  $d_j$
- Similarity is measured by

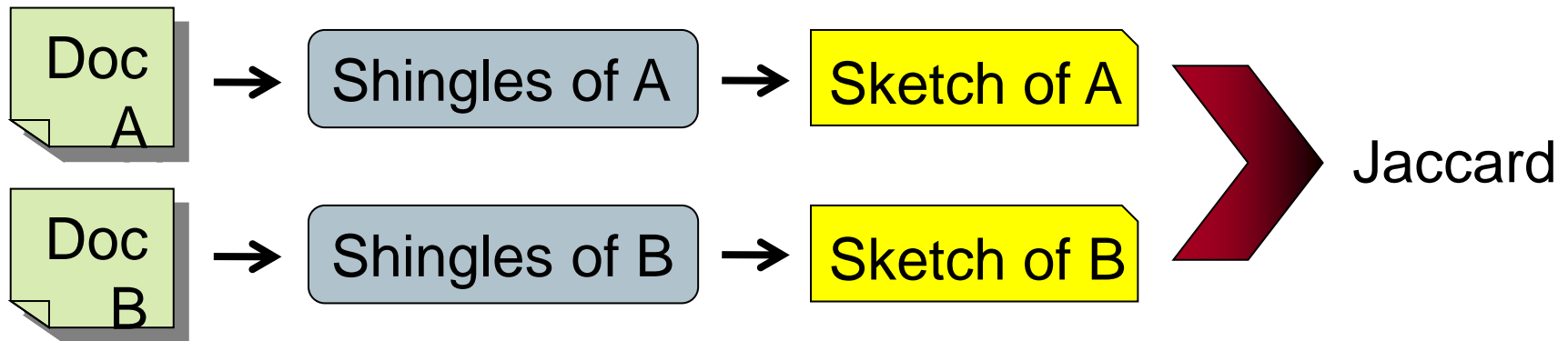
$$\text{Sim}(d_i, d_j) = \text{Jaccard}(S(d_i), S(d_j)) = \frac{|S(d_i) \cap S(d_j)|}{|S(d_i) \cup S(d_j)|}$$

- We have to compute Jaccard coefficients **for all pairs of web documents**.
- Computing exact set intersection of shingles between all pairs of web documents is **expensive** and **intractable**. 가

# Computing Similarity

---

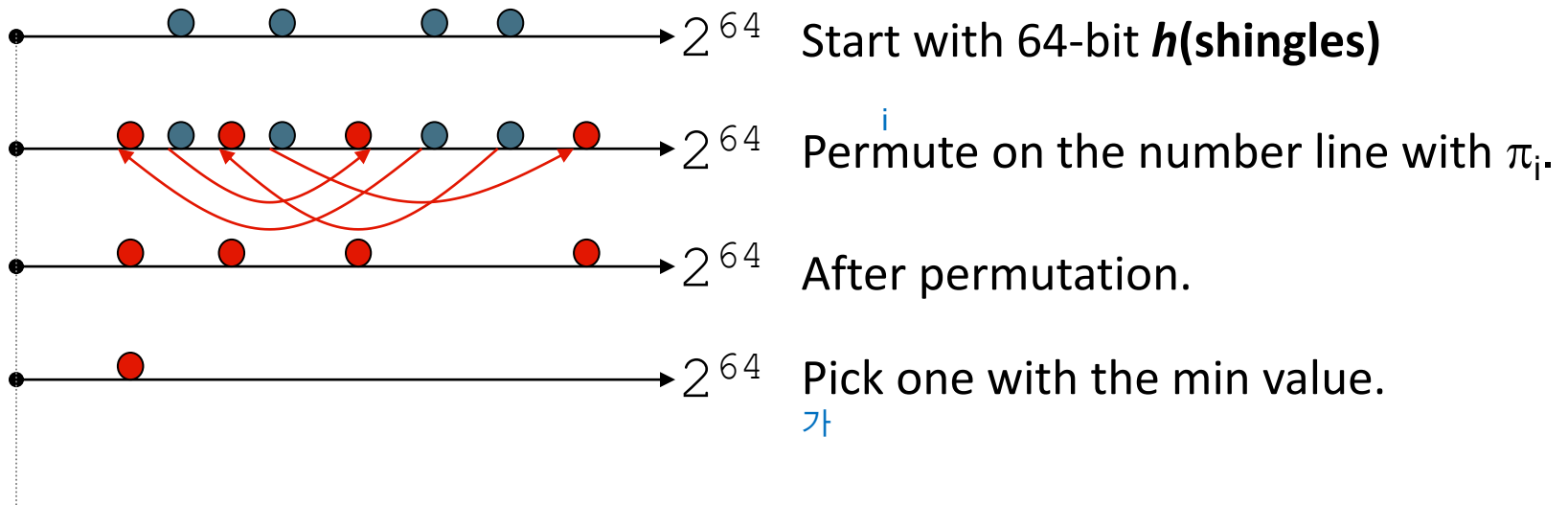
- **Approximate** using a cleverly chosen subset of shingles for each document.
- This subset is called a **sketch** of the document.
- Calculate similarity based on the short sketch.



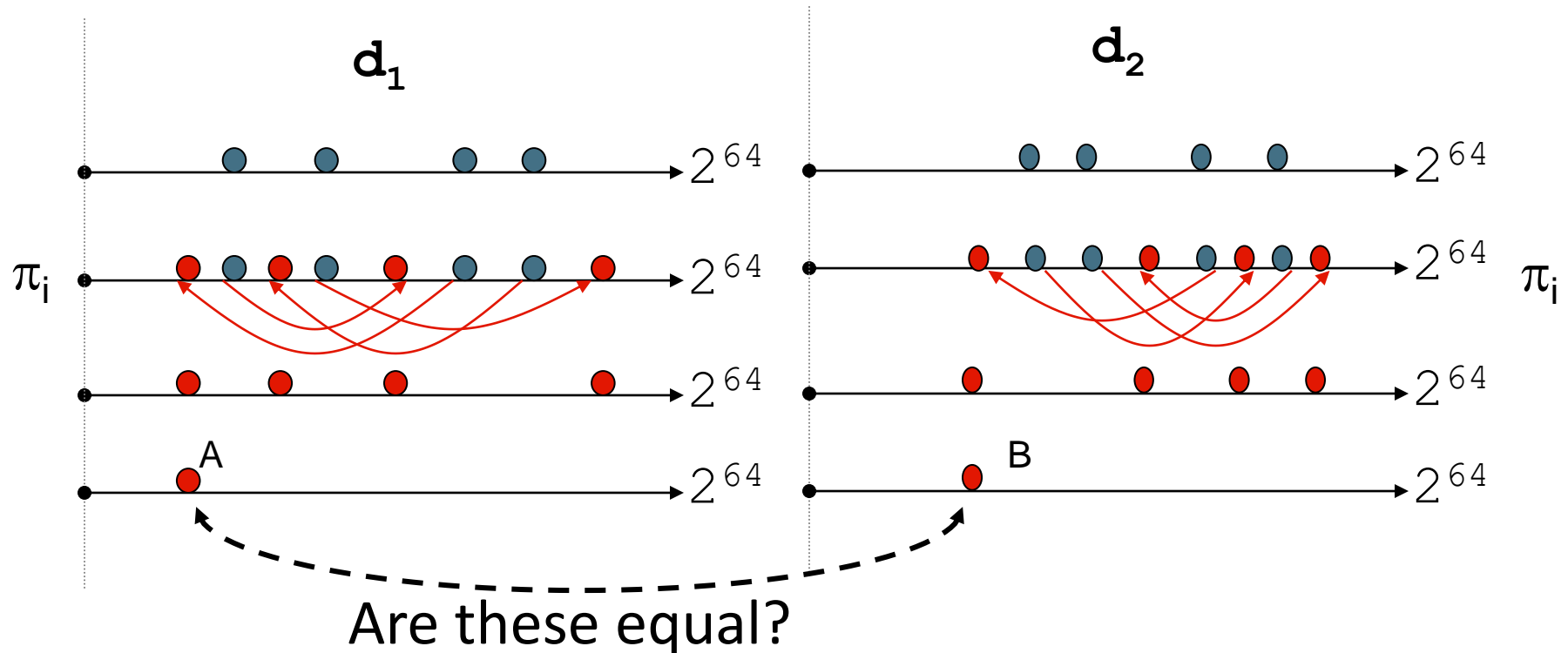
# Sketch of a Document <sup>h:</sup>

- For a document  $d$ , a  $sketch_d$  is given as follows:
  - Let  $h$  be a hash function that maps all shingles in the universe to  $[0..2^m]$ .

$h(\text{a\_rose\_is\_a}) = 921037246$
  - Let  $\pi_i$  be a random permutation on  $[0..2^m]$ .
  - Pick  $\text{MIN } \{\pi_i(h(s))\}$  over all shingles  $s$  for  $d$ . ( $\Rightarrow sketch_d$ )  
It is much like a *card shuffling* and choosing the first one.



# Test if $(Sketch_{d_1} == Sketch_{d_2})$



It is much like a *card shuffling* for both decks and choosing the first one from each deck, and seeing if they are the same one.



# Sketch of a Document

---

- Let  $S(d1) = \{\text{"apple", "banana", "berry"}\}$   
 $S(d2) = \{\text{"cherry", "apple", "lemon", "berry"}\}$
- After hashing

S(d1)	S(d2)
	cherry
berry	berry
	lemon
apple	apple
banana	

# Sketch of a Document

- After Permutation

S(d1)	S(d2)
	cherry
berry	berry
	lemon
apple	apple
banana	



S(d1)	S(d2)
	lemon
	cherry
apple	apple
berry	berry
banana	

**Sketch**

S(d1)	S(d2)
berry	berry
banana	
	cherry
apple	apple
	lemon

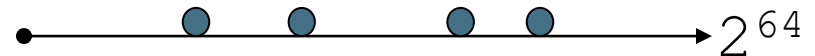
# Sketch of a Document

- Let  $S(d1) = \{\text{"red", "blue", "purple", "green", "yellow"}\}$   
 $S(d2) = \{\text{"blue", "orange", "purple", "violet"}\}$

- After hashing

$$S(d1) = \{1, 2, 4, 5, 7\}$$

$$S(d2) = \{2, 3, 4, 6\}$$



- Assume the following permutations

$$\pi_1 = \{(1 \rightarrow 7), (2 \rightarrow 2), (3 \rightarrow 4), (4 \rightarrow 1), (5 \rightarrow 3), (6 \rightarrow 6), (7 \rightarrow 5)\}$$

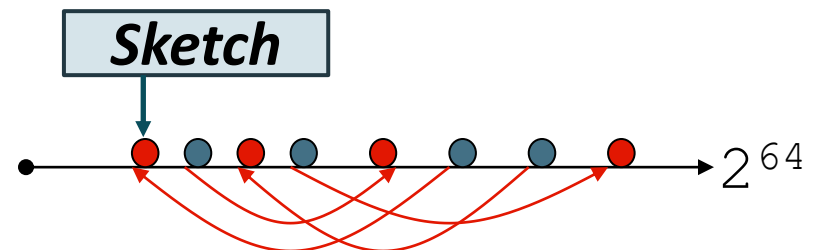
$$\pi_2 = \{(1 \rightarrow 4), (2 \rightarrow 6), (3 \rightarrow 5), (4 \rightarrow 2), (5 \rightarrow 1), (6 \rightarrow 7), (7 \rightarrow 3)\}$$

$$\pi_3 = \{(1 \rightarrow 3), (2 \rightarrow 4), (3 \rightarrow 1), (4 \rightarrow 5), (5 \rightarrow 2), (6 \rightarrow 7), (7 \rightarrow 6)\}$$

- Apply  $\pi_1$  to  $S(d1)$  and  $S(d2)$

$$\pi_1(S(d1)) = \{7, 2, \mathbf{1}, 3, 5\}$$

$$\pi_1(S(d2)) = \{2, 4, \mathbf{1}, 6\}$$



# Sketch of a Document

- After hashing

$$S(d1) = \{1, 2, 4, 5, 7\}, \quad S(d2) = \{2, 3, 4, 6\}$$

- Assume the following permutations

$$\pi_1 = \{(1 \rightarrow 7), (2 \rightarrow 2), (3 \rightarrow 4), (4 \rightarrow 1), (5 \rightarrow 3), (6 \rightarrow 6), (7 \rightarrow 5)\}$$

$$\pi_2 = \{(1 \rightarrow 4), (2 \rightarrow 6), (3 \rightarrow 5), (4 \rightarrow 2), (5 \rightarrow 1), (6 \rightarrow 7), (7 \rightarrow 3)\}$$

$$\pi_3 = \{(1 \rightarrow 3), (2 \rightarrow 4), (3 \rightarrow 1), (4 \rightarrow 5), (5 \rightarrow 2), (6 \rightarrow 7), (7 \rightarrow 6)\}$$

- Apply  $\pi_1, \pi_2, \pi_3$  to  $S(d1)$  and  $S(d2)$

$$\pi_1(S(d1)) = \{\mathbf{1}, 2, 3, 5, 7\}, \quad \pi_1(S(d2)) = \{\mathbf{1}, 2, 4, 6\}$$

$$\pi_2(S(d1)) = \{\mathbf{1}, 2, 3, 4, 6\}, \quad \pi_2(S(d2)) = \{\mathbf{2}, 5, 6, 7\}$$

$$\pi_3(S(d1)) = \{\mathbf{2}, 3, 4, 5, 6\}, \quad \pi_3(S(d2)) = \{\mathbf{1}, 4, 5, 7\}$$

permutation 3 : 1 가 1/3  
가  
1 1

# Computation of $Jaccard(S(d_1), S(d_2))$

---

- Theorem

$$Sim(d_1, d_2) = Jaccard(S(d_1), S(d_2)) = P(\text{sketch}_{d_1}, \text{sketch}_{d_2})$$

- Approximate  $Jaccard(S(d_1), S(d_2))$

- Create a **sketch vector**  $sketch_d[]$  (of size  $< 200$ ) for each document  $d_1$  and  $d_2$

→ Do 200 random permutations:  $\pi_1, \pi_2, \dots, \pi_{200}$

- If  $(sketch_{d_1}[i] == sketch_{d_2}[i])$  more than  $t$  times, then they are **near duplicates**.
- $t$  is determined empirically.

# Computation of $Jaccard(S(d_1), S(d_2))$

---

- Is sketch efficient?
  - Assume the  $S(d_1)$  and  $S(d_2)$  are shingles of two documents
  - To compute Jaccard of  $S(d_1)$  and  $S(d_2)$ , we have to know  $S(d_1) \cap S(d_2)$  which requires  $|S(d_1)| \times |S(d_2)|$  operations
  - Now using the sketch, the number of operation reduces  $|S(d_1)| + |S(d_2)|$  for each permutation. As we repeat the permutation  $P$  times, total operation required is:  
 $(|S(d_1)| + |S(d_2)|) \times P \ll |S(d_1)| \times |S(d_2)|$

# Computation of $Jaccard(S(d_1), S(d_2))$

---

- Is sketch efficient?
  - Assume that there are  $N$  documents in the collection
  - To compute Jaccard of  $S(d_i)$  and  $S(d_j)$ , we have to know  $S(d_i) \cap S(d_j)$  which requires  $|S(d_i)| \times |S(d_j)|$  operations
  - This should be repeat for all the pairs in the collection
  - Once a **sketch vector**  $sketch_d[]$  is created for each document  $d$  in the collection
  - Now using the sketch vectors, It is just 200 iterations of comparison to compute  $P(sketch_{d_i}, sketch_{d_j})$  for a pair of  $d_i$  and  $d_j$  in the collection

# Near Duplicate Documents

---

- Now we have an extremely efficient method for estimating a Jaccard coefficient for a pair of documents.
- But we still have to estimate  $N^2$  coefficients where  $N$  is the number of web pages.
  - Still slow!!

가

가