# Information Retrieval – Part 2

Deepak Kumar

# Borges' *Library of Babel*

"…each book contains four hundred ten pages; each page, forty lines; each line, approximately eighty black letters. There are also letters on the front cover of each book; these letters neither indicate nor prefigure what the pages inside will say."
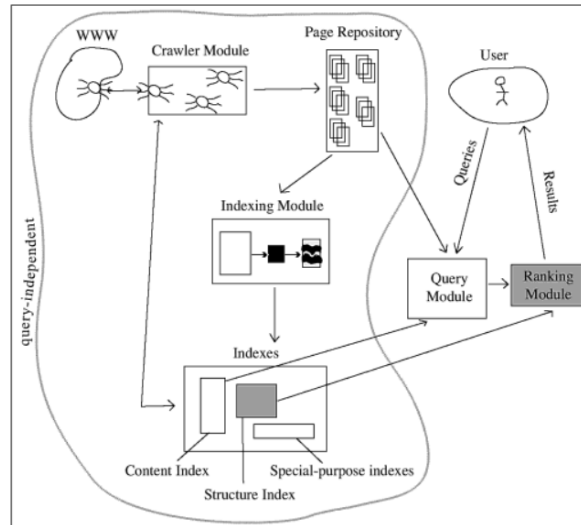
Q: How many books are in the library?
Q. How would you find what you're looking for?

# Elements of a Search Engine

# Web Information Retrieval

- Search Engines
- Queries
  phrase queries
  structure queries (NEAR, intitle:, …)
- Matching
- Inverted Index
  page number
  location
- Ranking & Relevance
- Metadata

# Web Information Retrieval

- Search Engines
- Queries
  phrase queries
  structure queries
- Matching
- Inverted Index
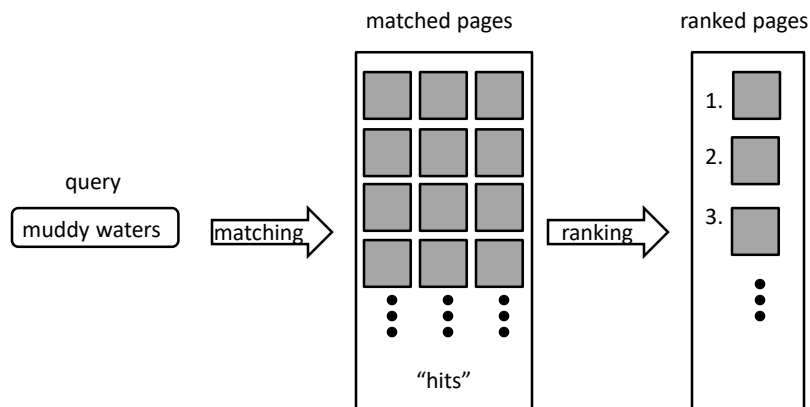  page number
  location
- Ranking & Relevance
- Metadata

**Efficient matching
is only one half the story.**

**The other grand challenge
is how to _rank_ the
matching pages**

5

# Matching & Ranking

matched pages                    ranked pages

query

muddy waters → matching → [hits] → ranking → 1. 2. 3.

"hits"

6

# Ranking & Relevance

1
By far the most common cause of malaria is being bitten by an infected mosquito, but there are also other ways to contract the disease.

2
Our cause was not helped by the poor health of the troops, many of whom were suffering from malaria and other tropical diseases.

query

malaria cause

| also | 1-19 | |
| ... | | |
| cause | 1-6 | 2-2 |
| ... | | |
| malaria | 1-8 | 2-19 |
| ... | | |
| whom | 2-15 | |

Nearness can resolve the ranking!

# Metadata
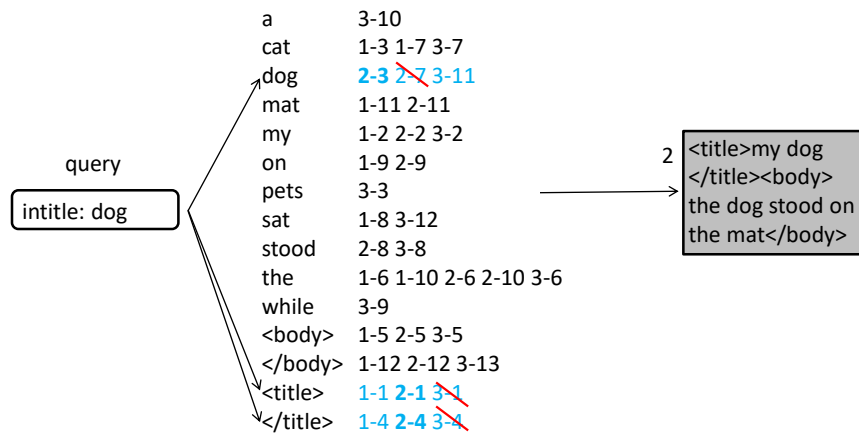
1
<title>my cat </title> <body> the cat sat on the mat </body>

2
<title>my dog </title><body> the dog stood on the mat</body>

3
<title>my pets </title><body>the cat stood while a dog sat

| a | 3-10 |
| cat | 1-3 1-7 3-7 |
| dog | 2-3 2-7 3-11 |
| mat | 1-11 2-11 |
| my | 1-2 2-2 3-2 |
| on | 1-9 2-9 |
| pets | 3-3 |
| sat | 1-8 3-12 |
| stood | 2-8 3-8 |
| the | 1-6 1-10 2-6 2-10 3-6 |
| while | 3-9 |
| <body> | 1-5 2-5 3-5 |
| </body> | 1-12 2-12 3-13 |
| <title> | 1-1 2-1 3-1 |
| </title> | 1-4 2-4 3-4 |

# Structure Queries

| | |
|---|---|
| a | 3-10 |
| cat | 1-3 1-7 3-7 |
| dog | **2-3** 2-7 3-11 |
| mat | 1-11 2-11 |
| my | 1-2 2-2 3-2 |
| on | 1-9 2-9 |
| pets | 3-3 |
| sat | 1-8 3-12 |
| stood | 2-8 3-8 |
| the | 1-6 1-10 2-6 2-10 3-6 |
| while | 3-9 |
| \<body\> | 1-5 2-5 3-5 |
| \</body\> | 1-12 2-12 3-13 |
| \<title\> | 1-1 **2-1** 3-1 |
| \</title\> | 1-4 **2-4** 3-4 |

query

intitle: dog

2  \<title\>my dog \</title\>\<body\> the dog stood on the mat\</body\>

# Exploiting Link Structure

- ***PageRank*** exploits the structure of the web:

  Use of Hyperlinks to
  - count # of incoming links
  - Identifying web authority

- Use the above in determining ranking & relevance.

# The Garage





**Garage at 232 Santa Margarita, Menlo Park, CA**

# Google 1.0 (1998)

2-proc Pentium II 300mhz, 512mb, five 9gb drives
2-proc Pentium II 300mhz, 512mb, four 9gb drives
4-proc PPC 604 333mhz, 512mb, eight 9gb drives
2-proc UltraSparc II 200mhz, 256mb, three 9gb
drives, six 4gb drives
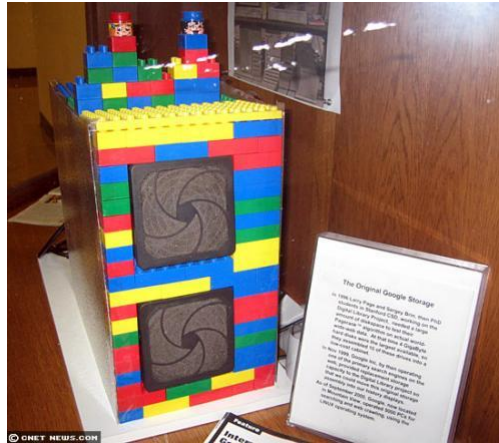Disk expansion, eight 9gb drives
Disk expansion, ten 9gb drives

That's a total of:
**1792 megabytes** of memory
**366 gigabytes** of disk storage
**2933 megahertz** in **10** CPUs

# The Disk Storage



13

# Google 1.0 (1998)



14

7

# Hyperlinks

**Ernie's Scrambled Eggs Recipe**
Mix four eggs in a bowl with a little salt and pepper, …

**Bert's Scrambled Eggs Recipe**
First melt a tablespoon of butter, …

Ernie's recipe is good.

I really enjoyed Bert's recipe.

Bert's recipe is wonderful!

Bert's recipe is fantastic!

15

# Hyperlinks: # Incoming Links

**ranked higher based on #incoming links**

**Ernie's Scrambled Eggs Recipe**
Mix four eggs in a bowl with a little salt and pepper, …

**Bert's Scrambled Eggs Recipe**
First melt a tablespoon of butter, …

Ernie's recipe is good.

I really enjoyed Bert's recipe.

Bert's recipe is wonderful!

Bert's recipe is fantastic!

16

8

# Hyperlinks: # Incoming Links

**ranked higher based on #incoming links**

**Ernie's Scrambled Eggs Recipe**
Mix four eggs in a bowl with a little salt and pepper, …

**Bert's Scrambled Eggs Recipe**
First melt a tablespoon of butter, …

Ernie's recipe is good.

I did not enjoy Bert's recipe.

Bert's recipe did not work.

Bert's recipe is unhealthy!

11/6/2019

17

# Hyperlinks: Authority

**Ernie's Scrambled Eggs Recipe**
Mix four eggs in a bowl with a little salt and pepper, …

**Bert's Scrambled Eggs Recipe**
First melt a tablespoon of butter, …

**Deepak Kumar's Home Page**
I tried Ernie's recipe once, and its not bad at all.

**Alton Brown's Home Page**
Bert's recipe is clearly one of the best.

11/6/2019

18

9

# Hyperlinks: Authority

**Ernie's Scrambled Eggs Recipe**
Mix four eggs in a bowl with a little salt and pepper, …  ②

**Bert's Scrambled Eggs Recipe**
First melt a tablespoon of butter, …  ⑩⓪

**Deepak Kumar's Home Page**
I tried Ernie's recipe once, and its not bad at all.  ②

**Alton Brown's Home Page**
Bert's recipe is clearly one of the best.  ⑩⓪

①  ①  ①  ①

…100 pages…

①  ①  ①  ①

①

①

# Cycles

A  B
C  D  E

# Computing Authority Scores

# Computing Authority Scores

# Computing Authority Scores

# Computing Authority Scores
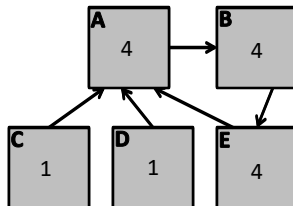
# Computing Authority Scores



# Computing Authority Scores

# Computing Authority Scores

# Computing Authority Scores

# Computing Authority Scores



and so on…stuck in an infinite loop….

# Sinks



**dangling node**

# Sinks

# The Random Surfer



restart probability = 15%

# The Random Surfer

**B**

| 109 |
|---|

| 74 |
|---|

| 55 |
|---|

| 55 |
|---|

| 32 |
|---|

| 58 |
|---|

| 16 |
|---|

**C**

| 13 |
|---|

| 46 |
|---|

| 17 |
|---|

| 21 |
|---|

| 15 |
|---|

| 101 |
|---|

**A**
| 118 |
|---|

| 126 |
|---|

| 144 | **D**
|---|

after 1000 page visits

11/6/2019

33

# The Random Surfer

**B**

| 10% |
|---|

| 7% |
|---|

| 4% |
|---|

| 4% |
|---|

| 4% |
|---|

| 5% |
|---|

| 2% |
|---|

**C**

| 2% |
|---|

| 4% |
|---|

| 2% |
|---|

| 2% |
|---|

| 2% |
|---|

| 10% |
|---|

**A**
| 13% |
|---|

| 14% |
|---|

| 15% | **D**
|---|

after 1 million page visits
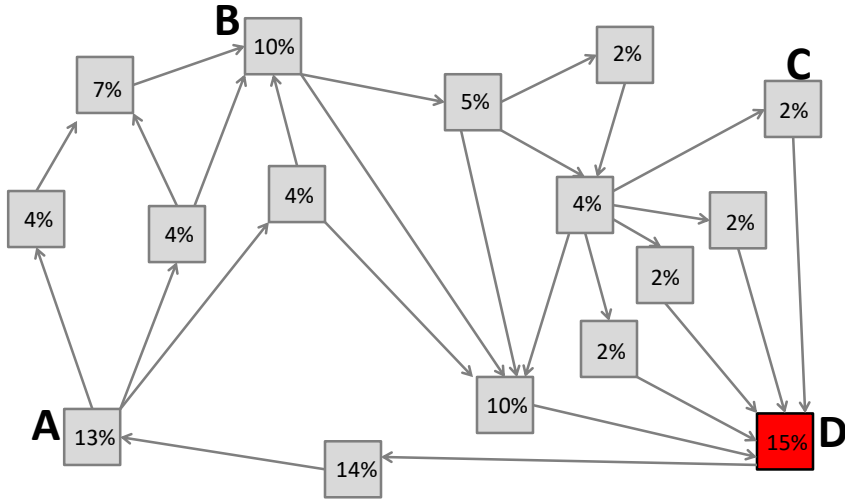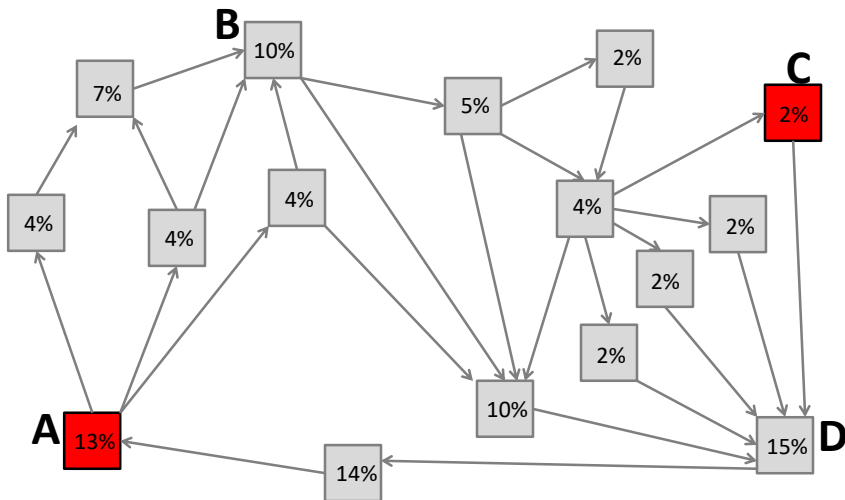
11/6/2019

34

17

The Random Surfer
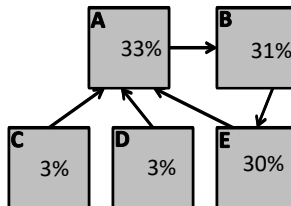
pages with many incoming links get high ranking



The Random Surfer

authoritative score

# The Random Surfer



#links+authoritative score

37

# The Random Surfer



**Ernie's Scrambled Eggs Recipe**
Mix four eggs in a bowl with a little salt and pepper, …
1%

**Bert's Scrambled Eggs Recipe**
First melt a tablespoon of butter, …
28%

0.4%

**Deepak Kumar's Home Page**
I tried Ernie's recipe once, and its not bad at all.
1%

0.4%

**Alton Brown's Home Page**
Bert's recipe is clearly one of the best.
32%

…100 pages…

0.4%

38

19

# The Random Surfer

---

# Formalizing PageRank

- Given a web page, $P_i$
- Set of pages pointing into $P_i$, $B_{P_i}$
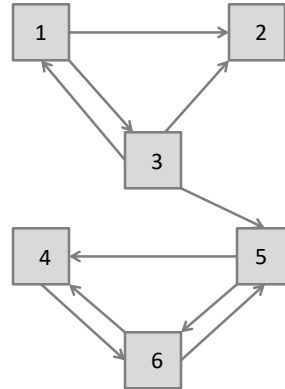- Number of outgoing links from page $P_j$, $|P_j|$
- PageRank of a page, $r(P_i)$

$$r(P_i) = \sum_{P_j \in B_{P_i}} \frac{r(P_j)}{|P_j|}$$

# Computing PageRank

- $r(P_1) = r(P_3)$

- But, $r(P_3)$ is unknown

- To start, assume all pages
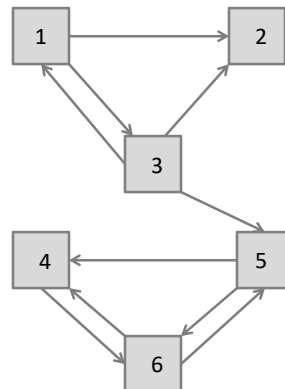  have rank $\frac{1}{n}$ $(n = 6)$

- $\therefore r(P_1) = \frac{1}{6}$

# Computing PageRank

$r_0(P_1) = 1/6$
$r_0(P_2) = 1/6$
$r_0(P_3) = 1/6$
$r_0(P_4) = 1/6$
$r_0(P_5) = 1/6$
$r_0(P_6) = 1/6$

# Computing PageRank

$r_1(P_1) = 1/18$
$r_1(P_2) = 5/36$
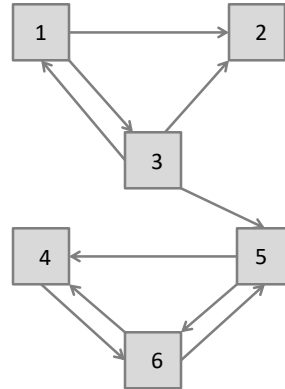$r_1(P_3) = 1/12$
$r_1(P_4) = 1/4$
$r_1(P_5) = 5/36$
$r_1(P_6) = 1/6$

$$r_{k+1}(P_i) = \sum_{P_j \in B_{P_i}} \frac{r_k(P_j)}{|P_j|}$$
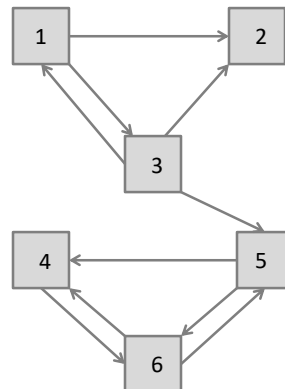
# Computing PageRank

$r_2(P_1) = 1/36$
$r_2(P_2) = 1/18$
$r_2(P_3) = 1/36$
$r_2(P_4) = 17/72$
$r_2(P_5) = 11/72$
$r_2(P_6) = 14/72$

$$r_{k+1}(P_i) = \sum_{P_j \in B_{P_i}} \frac{r_k(P_j)}{|P_j|}$$

# Computing PageRank

$r_2(P_1) = 1/36$      5
$r_2(P_2) = 1/18$      4
$r_2(P_3) = 1/36$      5
$r_2(P_4) = 17/72$     1
$r_2(P_5) = 11/72$     3
$r_2(P_6) = 14/72$     2

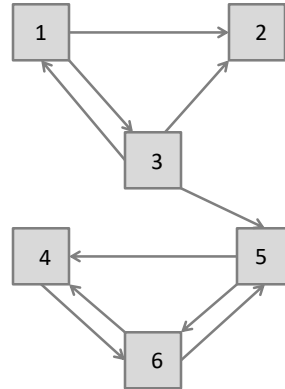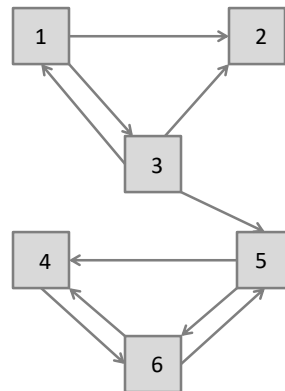$$r_{k+1}(P_i) = \sum_{P_j \in B_{P_i}} \frac{r_k(P_j)}{|P_j|}$$

# Matrix Representation

- Adjacency Matrix

$$A = \begin{array}{c} \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{array} \begin{array}{cccccc} 1 & 2 & 3 & 4 & 5 & 6 \\ \begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix} \end{array}$$
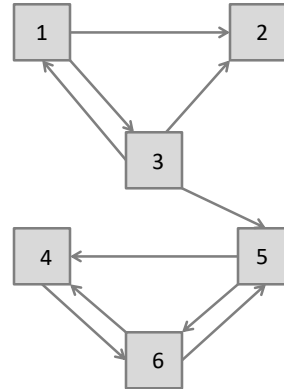
# Matrix Representation

- Hyperlink Matrix, $H$

$$\begin{array}{c} \phantom{1} \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{array} \begin{array}{cccccc} 1 & 2 & 3 & 4 & 5 & 6 \\ \left[\begin{array}{cccccc} 0 & 1/2 & 1/2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 1/3 & 1/3 & 0 & 0 & 1/3 & 0 \\ 0 & 0 & 0 & 0 & 1/2 & 1/2 \\ 0 & 0 & 0 & 1/2 & 0 & 1/2 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{array}\right] \end{array}$$



- $\pi_{k+1}{}^T = \pi_k{}^T H$

 where $\pi_k{}^T$ is the $k^{th}$ PageRank vector

# The PageRank Equation

- $\pi^T = \pi^T(\alpha S + (1 - \alpha)E)$

where
$S$ is the stochastic $H$ matrix
$E$ is the teleportation matrix
$\alpha$ is the scaling parameter

- Certain stochastic conditions apply!

## Google Data Center

## References

- *Google's PageRank and Beyond*, Amy N. Langville and Carl D. Meyer, Princeton University Press, 2006.
- *Nine Algorithms That Changed The Future*, John MacCormick, Princeton University Press, 2012.
- *The Unimaginable Mathematics of Borges' Library of Babel*, William G. Bloch, Oxford University Press, 2008.