



Information Retrieval Algorithms

National Tsing Hua University
2018, Fall Semester

Outline

■ Information retrieval

- Boolean Retrieval
- Ranked Retrieval
- Inverted indexing in MapReduce

■ Graphic Problem

- Parallel breadth-first search
- PageRank

Information retrieval (IR)

■ Information retrieval (IR)

- the activity of obtaining information resources relevant to an information need from a collection of information resources
- Focus on textual information (= text/document retrieval)
- Other possibilities include image, video, music, ...

■ What do we find?

- Generically, “documents”
- Even though we may be referring to web pages, PDFs, PowerPoint slides, paragraphs, etc.

The Central Problem in Search



Concepts



Query Terms
“tragic love story”



Author



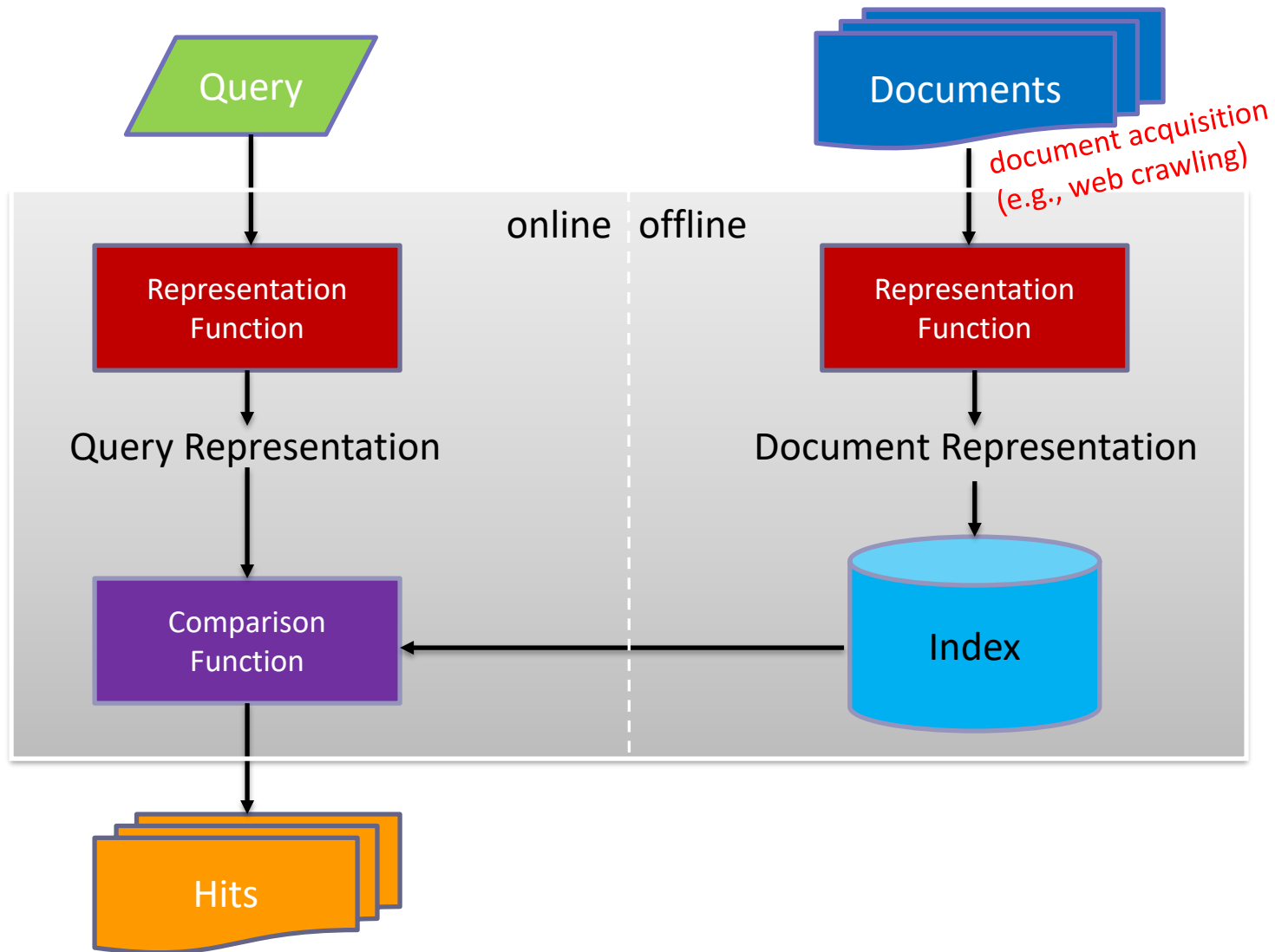
Concepts



Document Terms
“fateful star-crossed romance”

Do these represent the same concepts?

Abstract IR Architecture



How do we represent text?

- Remember: computers don't "understand" anything!
- "Bag of words"
 - Treat all the words in a document as **index** terms
 - Assign a "**weight**" to each term based on "**importance**" (or, in simplest case, presence/absence of word)
 - Disregard order, structure, meaning, etc. of the words
 - Simple, yet effective!
- Assumptions
 - Term occurrence is independent
 - Document relevance is independent
 - "Words" are well-defined

What's a word?

天主教教宗若望保祿二世因感冒再度住進醫院。這是他今年第二度因同樣的病因住院。

وقال مارك ريجيف - الناطق باسم
الخارجية الإسرائيلية - إن شارون قبل
الدعوة وسيقوم للمرة الأولى بزيارة
تونس، التي كانت لفترة طويلة المقر
الرسمي لمنظمة التحرير الفلسطينية بعد خروجها من لبنان عام 1982.

Выступая в Мещанском суде Москвы экс-глава ЮКОСа заявил
не совершал ничего противозаконного, в чем обвиняет его
генпрокуратура России.

भारत सरकार ने आर्थिक सर्वेक्षण में वित्तीय वर्ष 2005-06 में सात
फ़ीसदी विकास दर हासिल करने का आकलन किया है और कर सुधार
पर ज़ोर दिया है

日米連合で台頭中国に対処...アーミテージ前副長官提言

조재영 기자= 서울시는 25일 이명박 시장이 '행정중심복합도시'
건설안에 대해 '군대라도 동원해 막고싶은 심정'이라고 말했다는
일부 언론의 보도를 부인했다.

Sample Document

McDonald's slims down spuds

Fast-food chain to reduce certain types of fat in its french fries with new cooking oil.

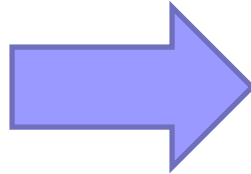
NEW YORK (CNN/Money) - McDonald's Corp. is cutting the amount of "bad" fat in its french fries nearly in half, the fast-food chain said Tuesday as it moves to make all its fried menu items healthier.

But does that mean the popular shoestring fries won't taste the same? The company says no. "It's a win-win for our customers because they are getting the same great french-fry taste along with an even healthier nutrition profile," said Mike Roberts, president of McDonald's USA.

But others are not so sure. McDonald's will not specifically discuss the kind of oil it plans to use, but at least one nutrition expert says playing with the formula could mean a different taste.

Shares of Oak Brook, Ill.-based McDonald's (MCD: down \$0.54 to \$23.22, Research, Estimates) were lower Tuesday afternoon. It was unclear Tuesday whether competitors Burger King and Wendy's International (WEN: down \$0.80 to \$34.91, Research, Estimates) would follow suit. Neither company could immediately be reached for comment.

...



"Bag of Words"

14 × McDonalds

12 × fat

11 × fries

8 × new

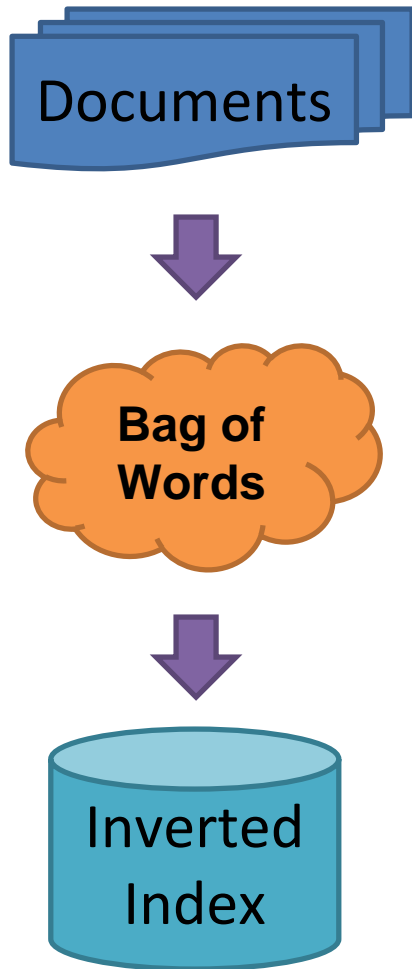
7 × french

6 × company, said, nutrition

5 × food, oil, percent, reduce,
taste, Tuesday

...

Counting Words...



1. Tokenization: Aren't → are not
2. Stopword removal: a, an, is, it, etc.
3. Normalization (equivalence classing of terms):
 - Hello → hello, windows → window, 你好 → hello

~~syntax~~, ~~semantics~~, ~~word knowledge~~, etc.

Outline

■ Information retrieval

- Boolean Retrieval
- Ranked Retrieval
- Inverted indexing in MapReduce

■ Graphic Problem

- Parallel breadth-first search
- PageRank

Boolean Retrieval

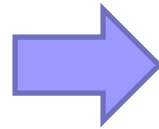
- Express queries as a **Boolean expression**
 - AND, OR, NOT
 - Can be arbitrarily nested
- Retrieval is based on the notion of sets
 - Any query divides the collection into **TWO sets**: retrieved, not-retrieved
 - Pure Boolean systems do **NOT** define an **ordering** of the results

Inverted Index: Boolean Retrieval

Indexed on words instead of documents

Doc 1 Doc 2 Doc 3 Doc 4
one fish, two fish red fish, blue fish cat in the hat green eggs and ham

	1	2	3	4
blue		1		
cat			1	
egg				1
fish	1	1		
green				1
ham				1
hat			1	
one	1			
red		1		
two	1			



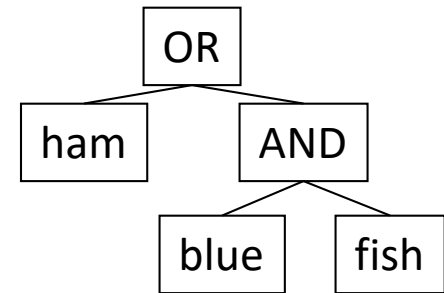
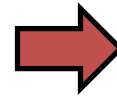
blue	→	2
cat	→	3
egg	→	4
fish	→	1 → 2
green	→	4
ham	→	4
hat	→	3
one	→	1
red	→	2
two	→	1

Boolean Retrieval

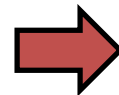
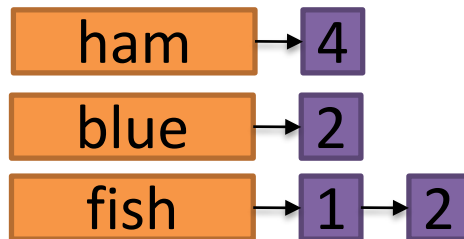
■ To execute a Boolean query:

- Build query syntax tree

(blue AND fish) OR ham



- For each clause, look up postings



{ {1,2} AND {2} } OR {4}



{2,4}

- Traverse postings and apply Boolean operator

■ Efficiency analysis

- Postings traversal is linear (assuming **sorted** postings)
- Start with shortest posting first

Strengths and Weaknesses

■ Strengths

- Precise, if you know the right strategies
- Precise, if you have an idea of what you're looking for
- Implementations are fast and efficient

■ Weaknesses

- Users must learn Boolean logic
- Boolean logic insufficient to capture the richness of language
- No control over the size of result set:
either too many hits or none
- **When do you stop reading?** All documents in the result set are considered “equally good”
- **What about partial matches?** Documents that “don't quite match” the query may be useful also

Outline

■ Information retrieval

- Boolean Retrieval
- Ranked Retrieval
- Inverted indexing in MapReduce

■ Graphic Problem

- Parallel breadth-first search
- PageRank

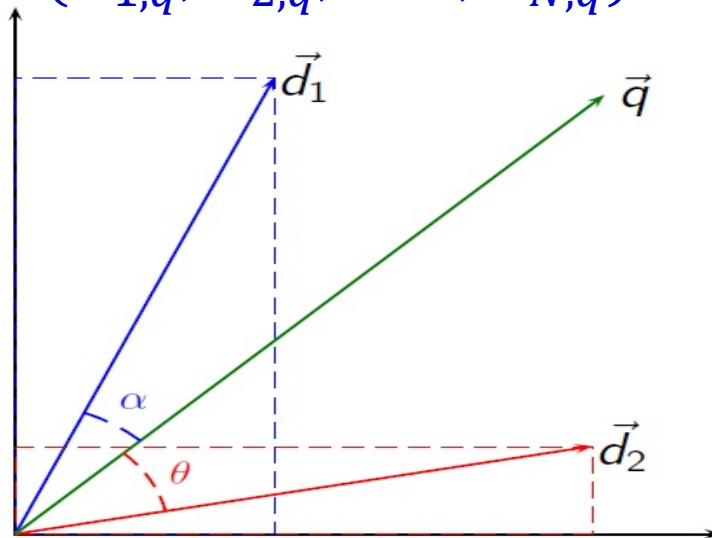
Ranked Retrieval

- Order documents by how likely they are to be relevant to the information need
 - Sort documents by relevance
 - Display sorted results
- User model
 - Present hits one screen at a time, best results first
 - At any point, users can decide to stop looking
- How do we estimate relevance?
 - Assume document is relevant if it has a lot of query **terms**
 - Represent query and document in **vector**
 - Relevance means the **closeness of vectors**

Vector Space Model

- Documents and queries are represented as vectors:
Each term is a dimension

- Document: $d_j = (w_{1,j}, w_{2,j}, \dots, w_{N,j})$
- Query: $q = (w_{1,q}, w_{2,q}, \dots, w_{N,q})$



retrieve documents based on how close the document is to the query
(i.e., similarity ~ “closeness”)

Similarity Metric

- Use “angle” between the vectors:

$$\text{sim}(d_j, q) = \cos(\theta) = \frac{\vec{d_j} \cdot \vec{q}}{|\vec{d_j}| \cdot |\vec{q}|}$$

$$\text{sim}(d_j, q) = \frac{\sum_{i=0}^n w_{i,j} w_{i,q}}{\sqrt{\sum_{i=0}^n (w_{i,j})^2} \sqrt{\sum_{i=0}^n (w_{i,q})^2}}$$

- Or, more generally, inner products:

$$\text{sim}(d_j, q) = \vec{d_j} \cdot \vec{q_k} = \sum_{i=0}^n w_{i,j} w_{i,q}$$

Term Weighting

- Term weights consist of two components
 - Local: how important is the term in this document?
 - Global: how important is the term in the collection?
- Here's the intuition:
 - Terms that appear often in a document should get high weights
 - Terms that appear in many documents should get low weights
- How do we capture this mathematically?
 - **Term frequency** (local)
 - Inverse **document frequency** (global)

TF.IDF Term Weighting

- Term Frequency-Inverse Document Frequency mode:
 - Terms appear often in a document should get high weights
 - Terms appear in many documents should get low weights

$$w_{i,j} = \text{tf}_{i,j} \cdot \log \frac{N}{df_i}$$

$w_{i,j}$ weight assigned to term i in document j

$\text{tf}_{i,j}$ frequency of occurrence of term i in document j

N number of documents in entire collection

df_i number of documents with term i

Inverted Index: TF.IDF

Doc 1

one fish, two fish

Doc 2

red fish, blue fish

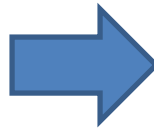
Doc 3

cat in the hat

Doc 4

green eggs and ham

		<i>tf</i>				<i>df</i>
		1	2	3	4	
blue			1			1
cat				1		1
egg					1	1
fish		2	2			2
green					1	1
ham					1	1
hat				1		1
one		1				1
red			1			1
two		1				1



		<i>df</i>		<i>tf</i>		
blue		1	2	1		
cat		1	3	1		
egg		1	4	1		
fish		2	1	2	2	2
green		1	4	1		
ham		1	4	1		
hat		1	3	1		
one		1	1	1		
red		1	2	1		
two		1	1	1		

Outline

■ Information retrieval

- Boolean Retrieval
- Ranked Retrieval
- Inverted indexing in MapReduce

■ Graphic Problem

- Parallel breadth-first search
- PageRank

MapReduce: Index Construction

- Map over all documents
 - Emit *term* as key, (*docno*, *tf*) as value
 - Emit other information as necessary (e.g., term position)
- Sort/shuffle:
 - group postings by term
- Reduce
 - Gather and sort the postings (e.g., by *docno* or *tf*)
 - Write postings to disk
- MapReduce does all the heavy lifting!

Inverted Indexing with MapReduce

Map

Doc 1	Doc 2	Doc 3
one fish, two fish	red fish, blue fish	cat in the hat
one	red	cat
two	blue	hat
fish	fish	

Shuffle and Sort: aggregate values by keys

Reduce

cat	<table><tr><td>3</td><td>1</td></tr></table>	3	1							
3	1									
fish	<table><tr><td>1</td><td>2</td></tr></table>	1	2	<table><tr><td>2</td><td>2</td></tr></table>	2	2	blue	<table><tr><td>2</td><td>1</td></tr></table>	2	1
1	2									
2	2									
2	1									
one	<table><tr><td>1</td><td>1</td></tr></table>	1	1		hat	<table><tr><td>3</td><td>1</td></tr></table>	3	1		
1	1									
3	1									
red	<table><tr><td>2</td><td>1</td></tr></table>	2	1		two	<table><tr><td>1</td><td>1</td></tr></table>	1	1		
2	1									
1	1									

Inverted Indexing: Pseudo-Code

```
1: class MAPPER
2:   procedure MAP(docid  $n$ , doc  $d$ )
3:      $H \leftarrow$  new ASSOCIATIVEARRAY
4:     for all term  $t \in$  doc  $d$  do
5:        $H\{t\} \leftarrow H\{t\} + 1$ 
6:     for all term  $t \in H$  do
7:       EMIT(term  $t$ , posting  $\langle n, H\{t\} \rangle$ )
```

$H\{t\}$: term frequency in a file

```
1: class REDUCER
2:   procedure REDUCE(term  $t$ , postings  $[\langle a_1, f_1 \rangle, \langle a_2, f_2 \rangle \dots]$ )
3:      $P \leftarrow$  new LIST
4:     for all posting  $\langle a, f \rangle \in$  postings  $[\langle a_1, f_1 \rangle, \langle a_2, f_2 \rangle \dots]$  do
5:       APPEND( $P, \langle a, f \rangle$ )
6:       SORT( $P$ )
7:       EMIT(term  $t$ , postings  $P$ )
```

Problem?

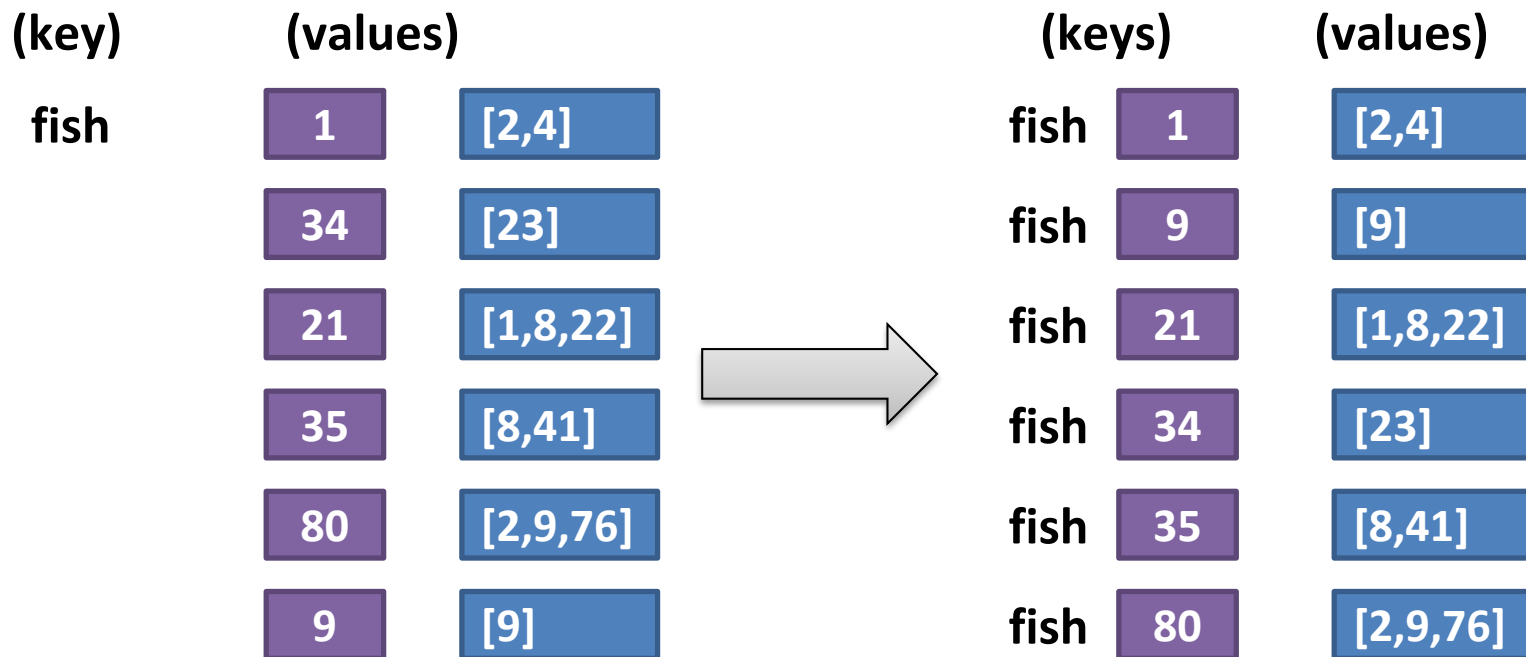
Scalability Bottleneck

- Initial ideal: terms as keys, postings as values
 - Reducers must buffer all postings associated with key (to sort)
 - What if we run **out of memory** to buffer postings?
- Think about 100million documents....
- Each posting can be long for storing additional info
 - Term position in a doc
 - Article title
 - Snippet

blue 2 1 1st paragraph, 453th word

Another Try...

- Use **<term><docid>** as the key
 - Let the framework do the sorting
 - Term frequency implicitly stored
 - **Directly write postings to disk!**



Retrieval with MapReduce?

- MapReduce is fundamentally **batch-oriented**
 - Optimized for throughput, not latency
 - Startup of mappers and reducers is expensive
- MapReduce is not suitable for **real-time** queries!
 - Initial a job takes time
 - Use separate infrastructure for retrieval...

Outline

■ Information retrieval

- Boolean Retrieval
- Ranked Retrieval
- Inverted indexing in MapReduce

■ Graphic Problem

- Parallel breadth-first search
- PageRank

Some Graph Problems

- Finding shortest paths
 - Routing Internet traffic and UPS trucks
- Finding minimum spanning trees
 - Telco laying down fiber
- Finding Max Flow
 - Airline scheduling
- Bipartite matching
 - Monster.com, Match.com
- And of course... PageRank

Graphs and MapReduce

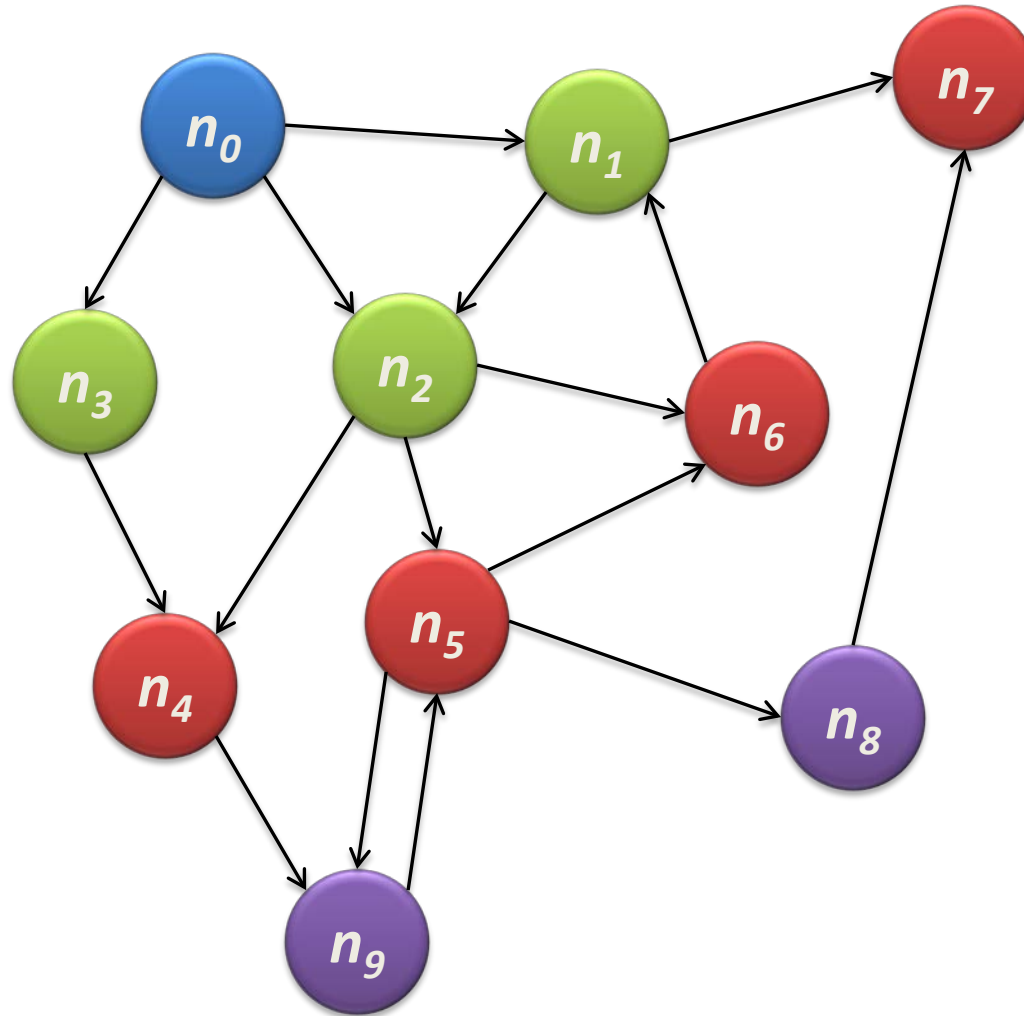
■ Graph algorithms typically involve:

- Performing computations at each node: based on node features, edge features, and local link structure
- Propagating computations: “traversing” the graph

■ Key questions:

- How do you represent graph data in MapReduce?
- How do you traverse a graph in MapReduce?

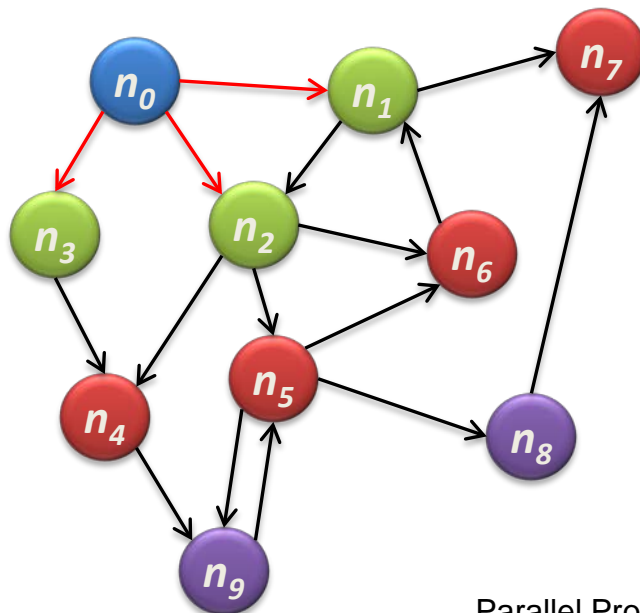
Visualizing Parallel BFS



Single Source Shortest Path

■ Data representation:

- Key: node n
- Value: d (distance from start),
adjacency list (list of nodes reachable from n)
- Initialization: for all nodes except for start node, $d = \infty$



n0,	0:n1:n2:n3
n1,	INF:n2:n7
n2,	INF:n4:n5:n6
n3,	INF:n4
n4,	INF:n9
n5,	INF:n6:n8:n9
n6,	INF:n1
n7,	INF:
n8,	INF:n7
n9,	INF:n5

Single Source Shortest Path

■ Mapper:

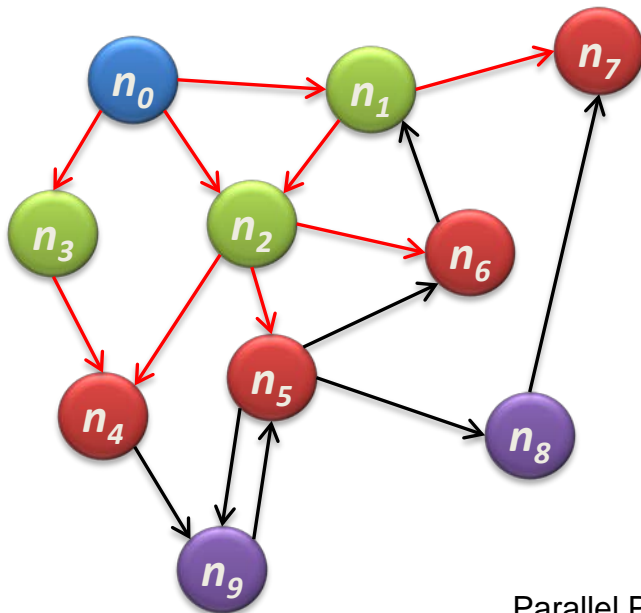
- $\forall m \in \text{adjacency list: emit } (m, d + 1)$

■ Sort/Shuffle

- Groups distances by **reachable nodes**

■ Reducer:

- **Selects minimum distance** path for each reachable node



n0,0:n1:n2:n3

n1,1:n2:n7

n2,1:n4:n5:n6

n3,1:n4

n4,INF:n9

n5,INF:n6:n8:n9

n6,INF:n1

n7,INF:

n8,INF:n7

n9,INF:n5

n1,1

n2,1

n3,1

n2,2

n7,2

n4,2

n5,2

n6,2

n4,2

n1,1

n2,1

n2,2

n3,1

n4,2

n4,2

n5,2

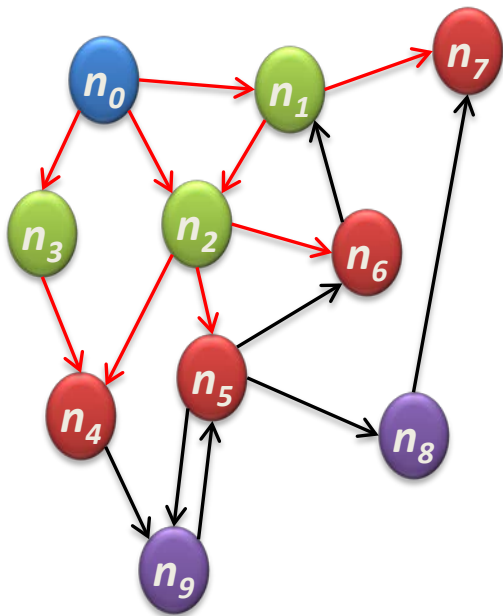
n6,2

n7,2

Multiple Iterations Needed

- Each MapReduce iteration advances the “known frontier” by one hop
 - Subsequent iterations include more and more reachable nodes as frontier expands
 - Multiple iterations are needed to explore entire graph
- Preserving graph structure:
 - Problem: Where did the adjacency list go?
 - Solution: mapper emits (n , adjacency list) as well

Start 2nd iter:



MAP

n0,0:n1:n2:n3
n1,1:n2:n7
n2,1:n4:n5:n6
n3,1:n4
n4,INF:n9
n5,INF:n6:n8:n9
n6,INF:n1
n7,INF:
n8,INF:n7
n9,INF:n5

n0,0:n1:n2:n3

n1,1:n2:n7
n1,1

n2,1:n4:n5:n6
n2,1
n2,2

n3,1:n4
n3,1

n4,INF:n9
n4,2
n4,2

n5,INF:n6:n8:n9
n5,2

n6,INF:n1
n6,2

n7,INF:
n7,2

n8,INF:n7

n9,INF:n5

REDUCE

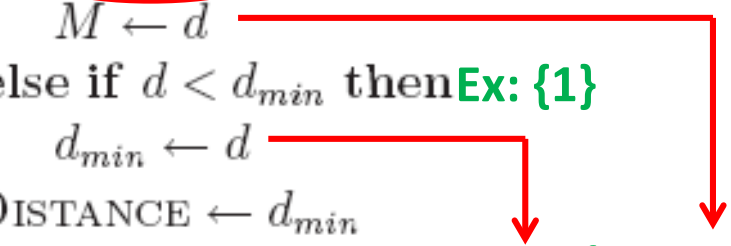
n0,0:n1:n2:n3
n1,1:n2:n7
n2,1:n4:n5:n6
n3,1:n4
n4,2:n9
n5,2:n6:n8:n9
n6,2:n1
n7,2:
n8,INF:n7
n9,INF:n5

Next
Iteration

BFS Pseudo-Code

```
1: class MAPPER
2:   method MAP(nid  $n$ , node  $N$ )
3:      $d \leftarrow N.DISTANCE$ 
4:     EMIT(nid  $n$ ,  $N$ ) Ex:  $\langle n0, \{0, n1:n2:n3\} \rangle$   $\triangleright$  Pass along graph structure
5:     for all nodeid  $m \in N.ADJACENCYLIST$  do
6:       EMIT(nid  $m$ ,  $d + 1$ )  $\triangleright$  Emit distances to reachable nodes

1: class REDUCER Ex:  $\langle n1, 1 \rangle$ ,  $\langle n2, 1 \rangle$ ,  $\langle n3, 1 \rangle$ 
2:   method REDUCE(nid  $m$ , [ $d_1, d_2, \dots$ ]) Ex:  $n1, \{3, n2:n4\}, \{1\}, \{4\}$ 
3:      $d_{min} \leftarrow \infty$ 
4:      $M \leftarrow \emptyset$ 
5:     for all  $d \in \text{counts } [d_1, d_2, \dots]$  do
6:       if IsNODE( $d$ ) then Ex:  $\{3, n2:n4\}$ 
7:          $M \leftarrow d$   $\triangleright$  Recover graph structure
8:       else if  $d < d_{min}$  then Ex:  $\{1\}$   $\triangleright$  Look for shorter distance
9:          $d_{min} \leftarrow d$ 
10:     $M.DISTANCE \leftarrow d_{min}$   $\triangleright$  Update shortest distance
11:    EMIT(nid  $m$ , node  $M$ ) Ex:  $\langle n1, \{1, n2:n4\} \rangle$ 
```



Stopping Criterion

■ Execution of an iterative MapReduce algorithm

- Typically, requires a **non-MapReduce driver program**, which submits a MapReduce job to iterate the algorithm until a **termination condition** has been met

■ How many iterations are needed

- Convince yourself: when a node is first “discovered”, we’ve found the shortest path
- Graph diameter (unit edge length)
- Number of nodes (weighted edge)
- Six degrees of separation: everyone on the planet is connected to everyone else by at most six steps

Outline

■ Information retrieval

- Boolean Retrieval
- Ranked Retrieval
- Inverted indexing in MapReduce

■ Graphic Problem

- Parallel breadth-first search
- PageRank

Random Walks Over the Web

■ Random surfer model:

- User starts at a random Web page
- User randomly clicks on links, surfing from page to page

■ PageRank

- Characterizes the amount of time spent on any given page
- Mathematically, a probability distribution over pages

■ PageRank captures notions of page importance

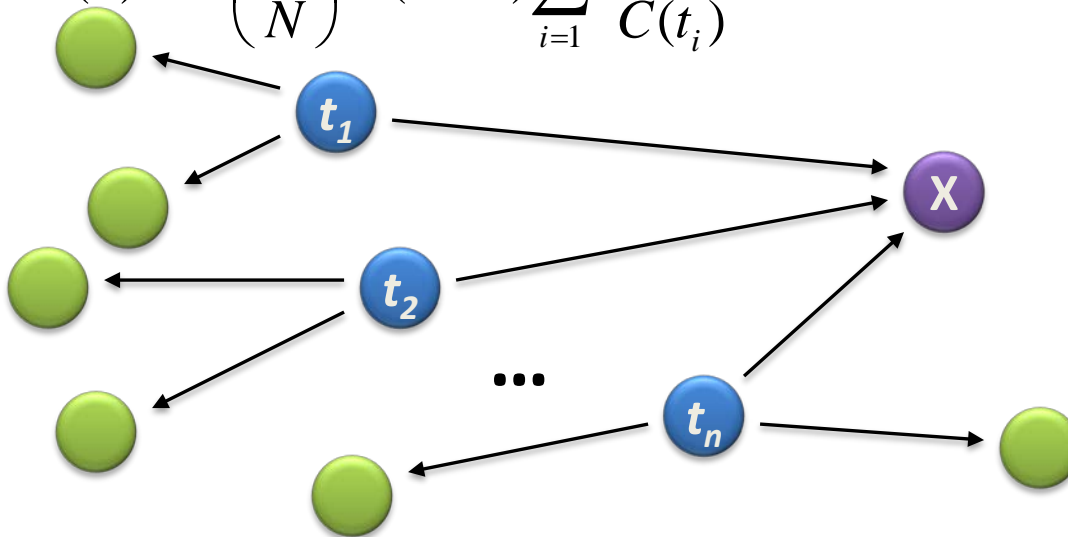
- Correspondence to human intuition?
- One of thousands of features used in web search
- Note: query-independent

PageRank: Defined

Given page x with inlinks $t_1 \dots t_n$, where

- $C(t)$ is the out-degree of t
- α is probability of **random jump**
- N is the total number of nodes in the graph

$$PR(x) = \alpha \left(\frac{1}{N} \right) + (1 - \alpha) \sum_{i=1}^n \frac{PR(t_i)}{C(t_i)}$$



Computing PageRank

■ Properties of PageRank

- Can be computed iteratively
- Effects at each iteration are local

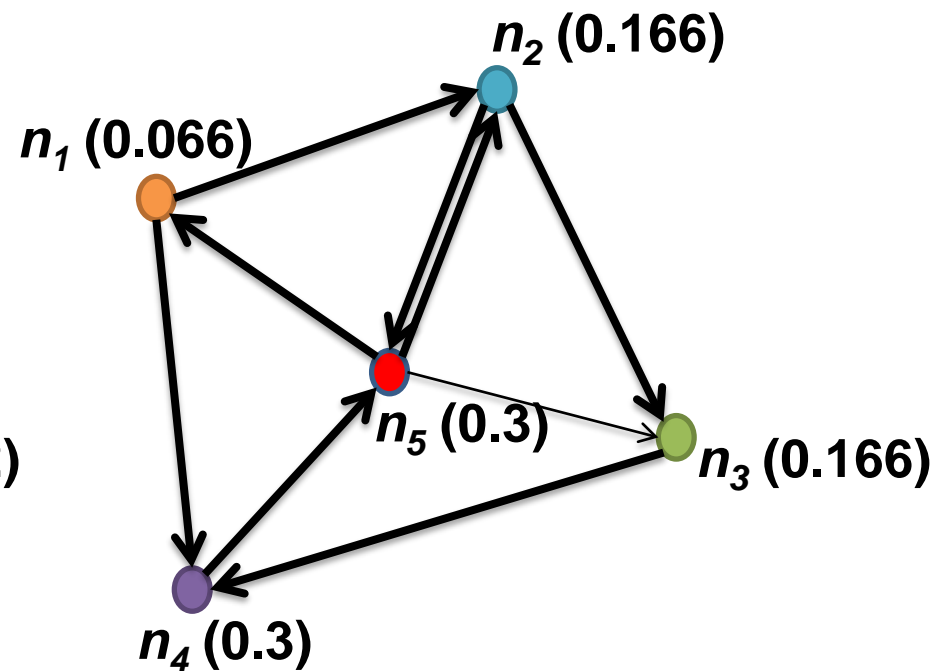
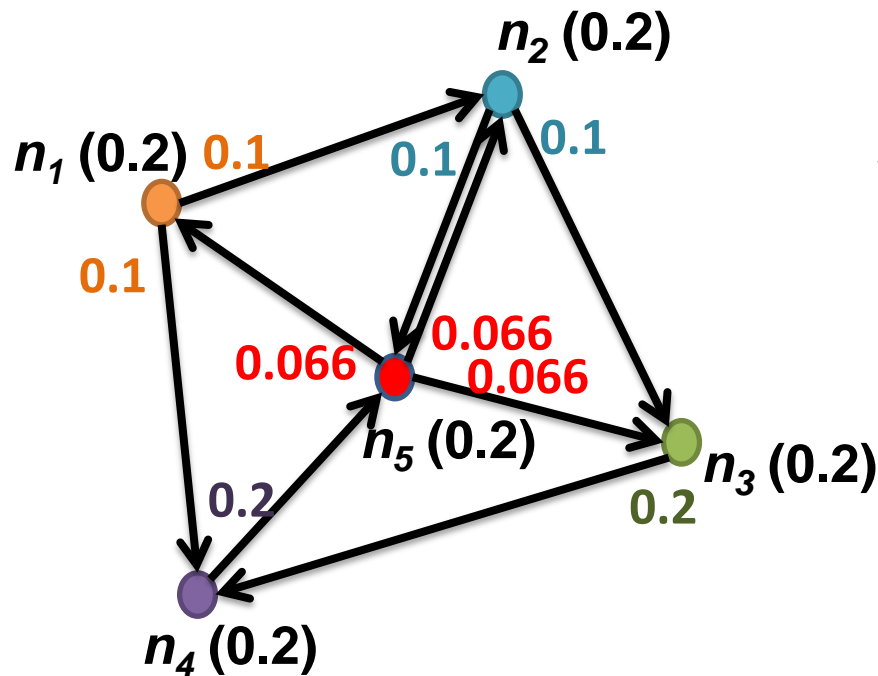
■ Sketch of algorithm:

- Start with seed PR_i values
- Each page distributes PR_i “credit” to all pages it links to
- Each target page adds up “credit” from multiple in-bound links to compute PR_{i+1}
- Iterate until values **converge**

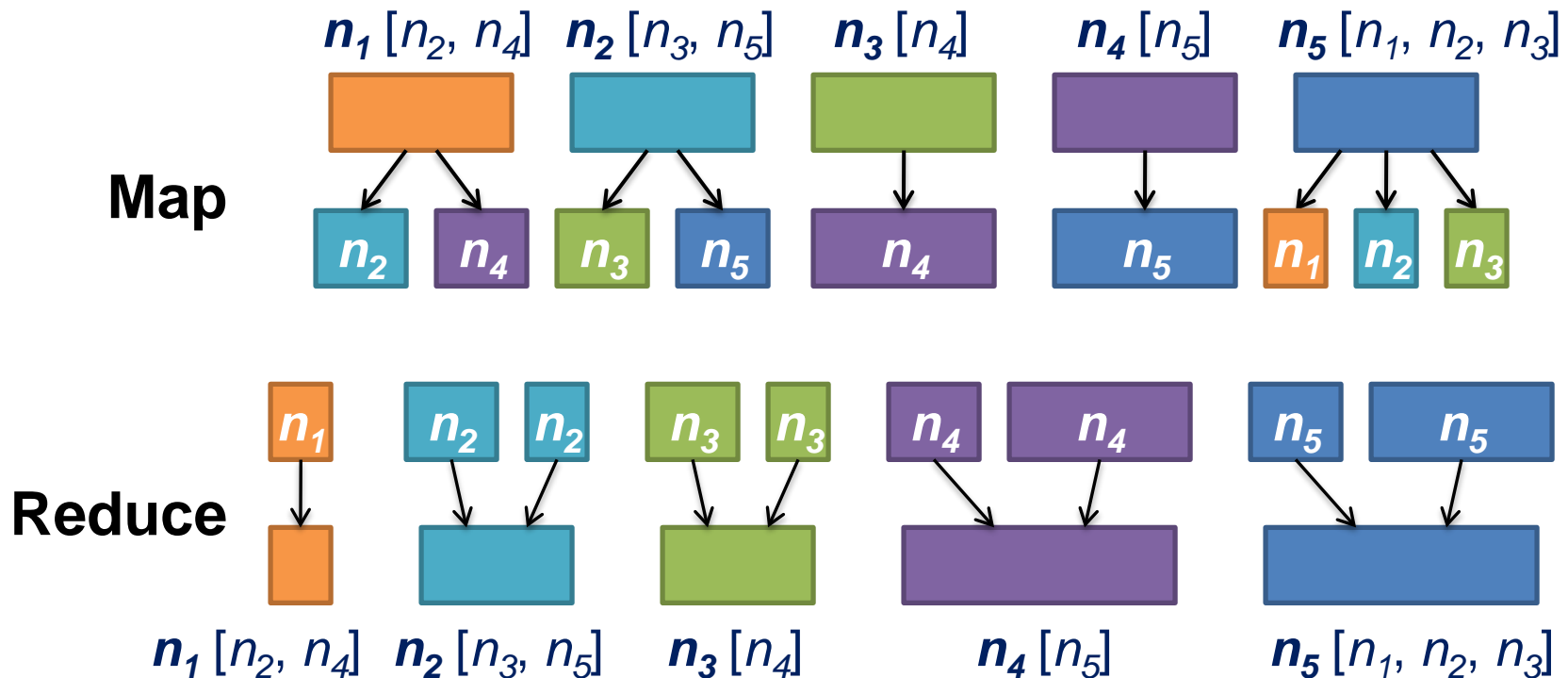
PageRank Example

- Simple case without random walk: $\alpha = 0$

$$PR(x) = \alpha \left(\frac{1}{N} \right) + (1 - \alpha) \sum_{i=1}^n \frac{PR(t_i)}{C(t_i)}$$



PageRank in MapReduce



PageRank Pseudo-Code

```
1: class MAPPER
2:   method MAP(nid  $n$ , node  $N$ )
3:      $p \leftarrow N.PAGERANK / |N.ADJACENCYLIST|$ 
4:     EMIT(nid  $n$ ,  $N$ )                                ▷ Pass along graph structure
5:     for all nodeid  $m \in N.ADJACENCYLIST$  do
6:       EMIT(nid  $m$ ,  $p$ )                                ▷ Pass PageRank mass to neighbors

1: class REDUCER
2:   method REDUCE(nid  $m$ , [ $p_1, p_2, \dots$ ])
3:      $M \leftarrow \emptyset$ 
4:     for all  $p \in$  counts [ $p_1, p_2, \dots$ ] do
5:       if ISNODE( $p$ ) then
6:          $M \leftarrow p$                                 ▷ Recover graph structure
7:       else
8:          $s \leftarrow s + p$                                 ▷ Sums incoming PageRank contributions
9:      $M.PAGERANK \leftarrow s$ 
10:    EMIT(nid  $m$ , node  $M$ )
```

PageRank Convergence

■ Alternative convergence criteria

- Iterate until PageRank values don't change
- Iterate until PageRank rankings don't change
- Fixed number of iterations

Reference

- http://hadoop.apache.org/common/docs/r1.0.3/mapred_tutorial.html
- <http://hadoop.apache.org/common/docs/r1.0.3/api/org/apache/hadoop/mapred/JobConf.html>
- <http://developer.yahoo.com/hadoop/tutorial/module5.html>
- **The PageRank Citation Ranking: Bringing Order to the Web.** Page, Lawrence and Brin, Sergey and Motwani, Rajeev and Winograd, Terry (1999) Technical Report. Stanford InfoLab.