

# THE STRUCTURE OF THE GOOGLE SEARCH ENGINE

**By Mohamed Ali**

# PRE GOOGLE

- Human Maintained Indices



- Subjective
- Expensive
- Cannot cover esoteric topics

- Automated Search Engines

- Relies chiefly on keyword matching
- Can be misled

# HOW BAD IS KEYWORD MATCHING

In 1997, Only one of the top four commercial search engines finds itself.

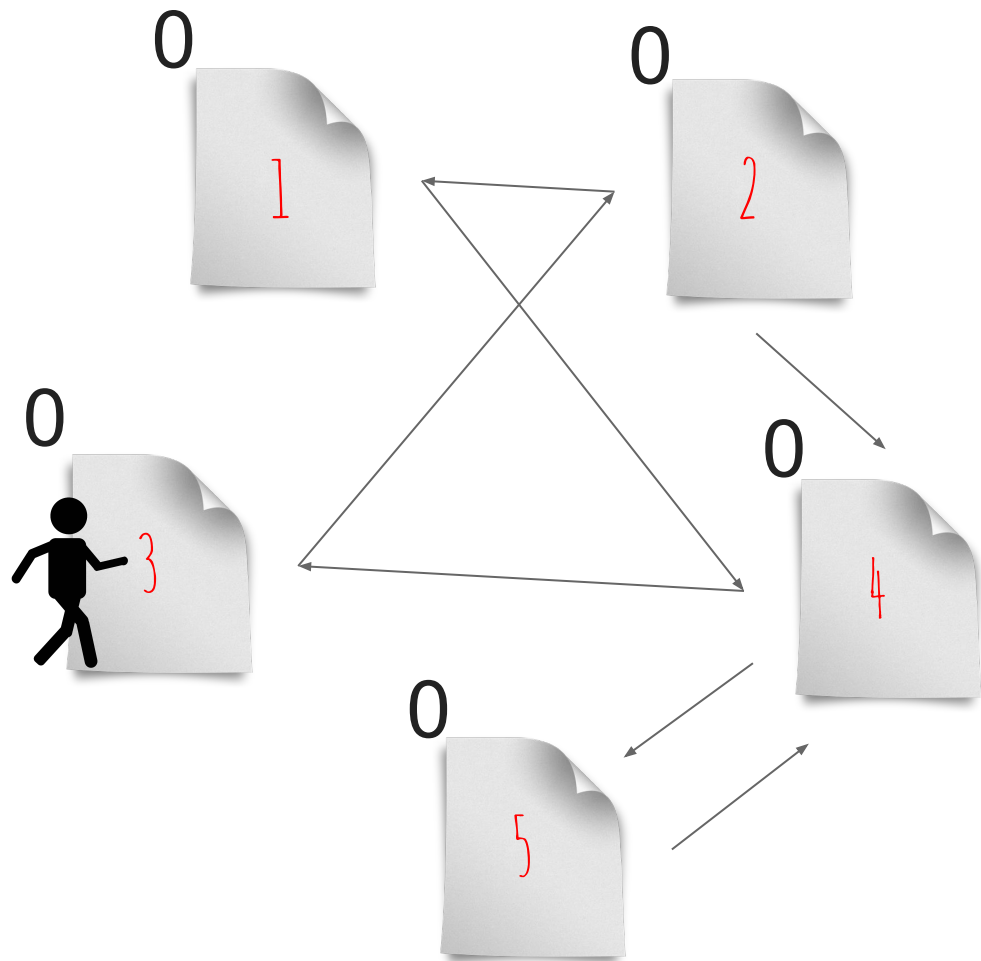
INTRODUCING  
GOOGLE

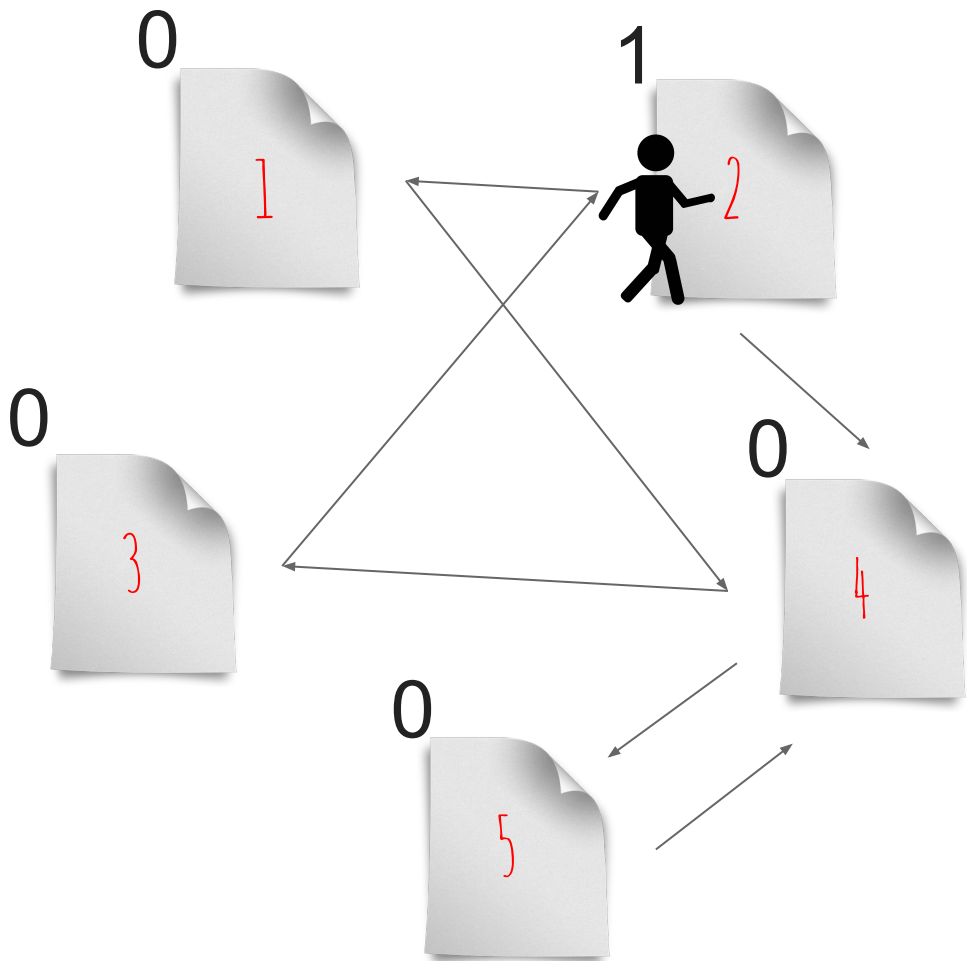
# SYSTEM FEATURES

- Making use of the link structure
  - Ordering the web
- Using anchor text
  - Sometimes more descriptive than the website's title
  - Allows for reaching non-text pages

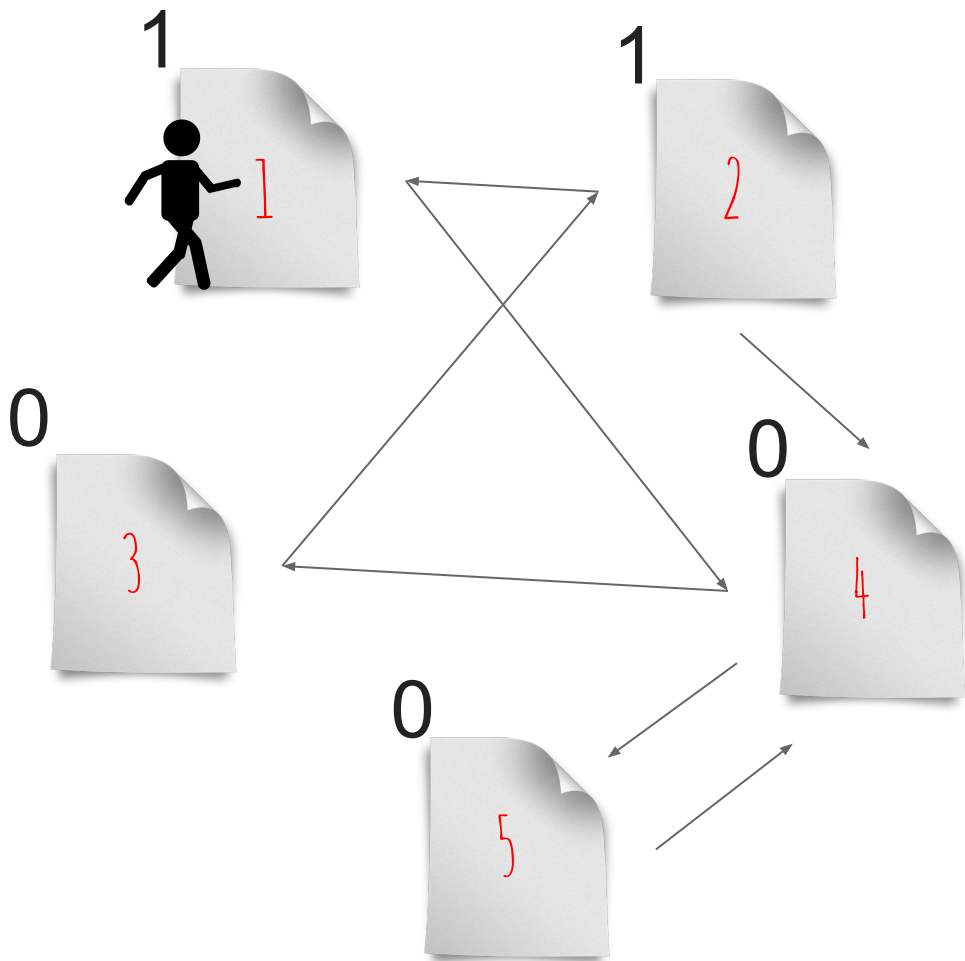
# CRAWLING AND INDEXING

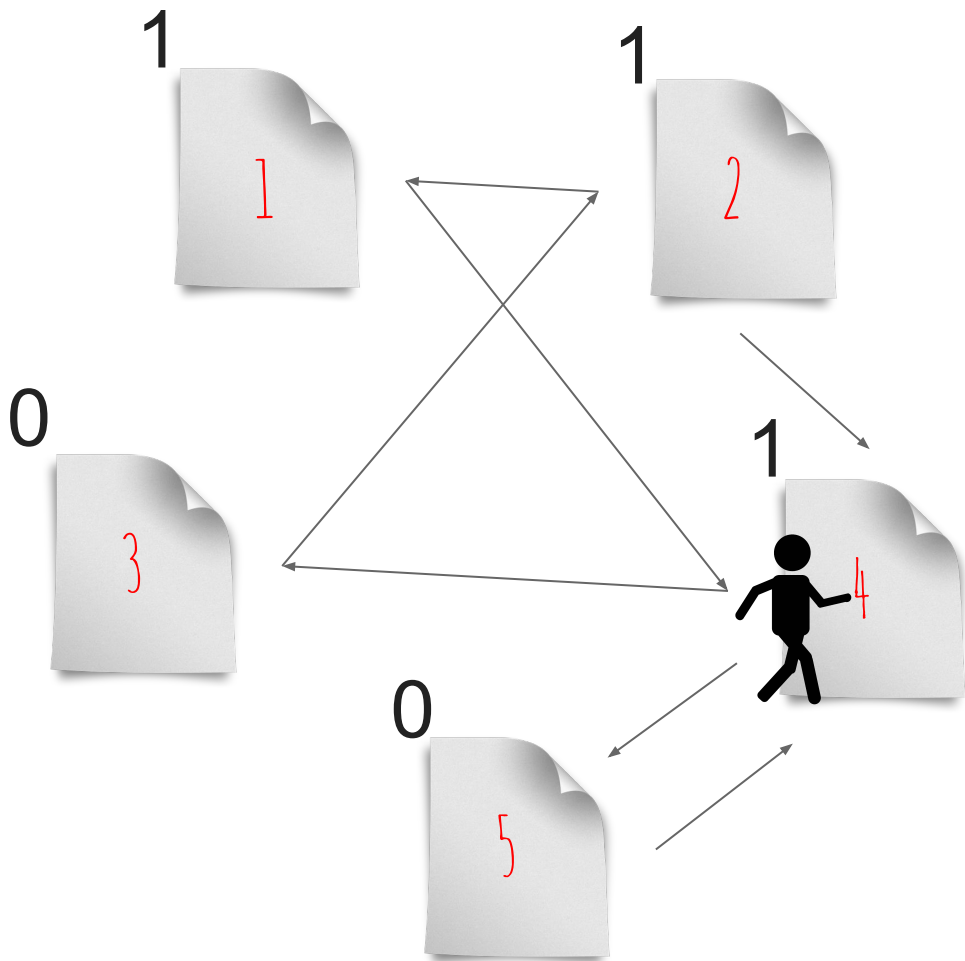
- The Random Surfer Model (PageRank)

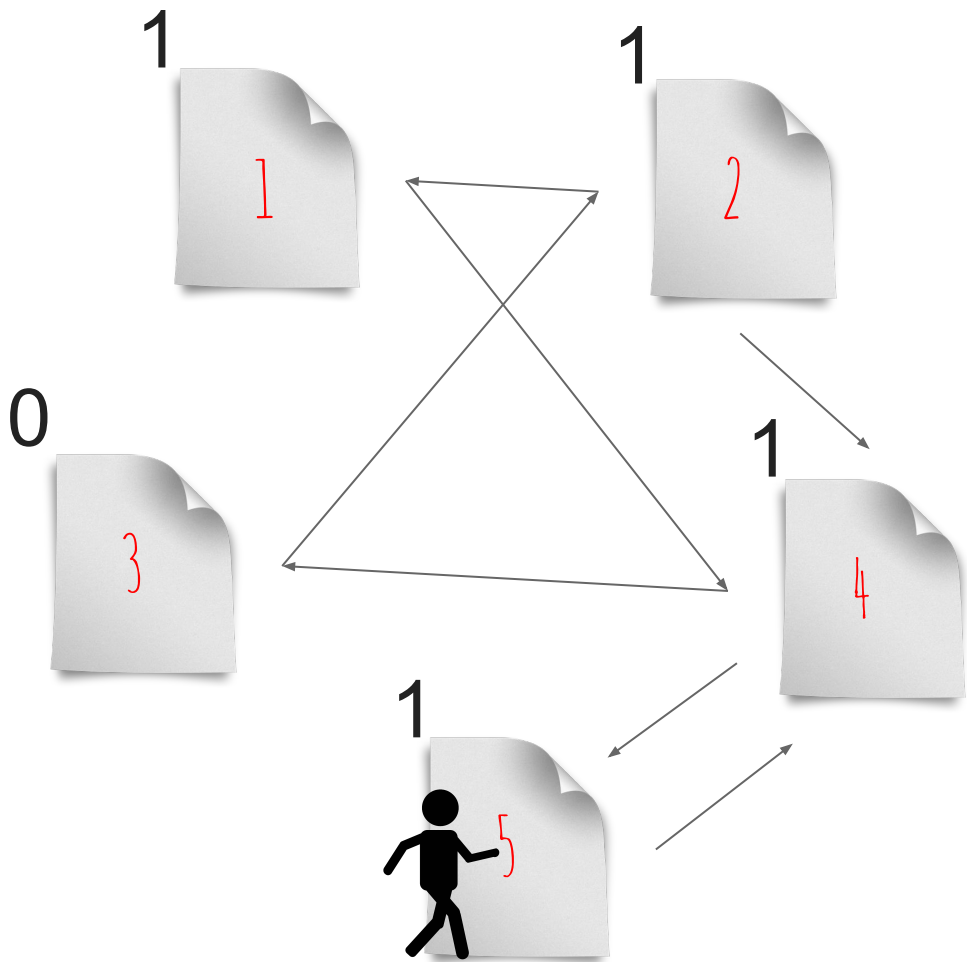


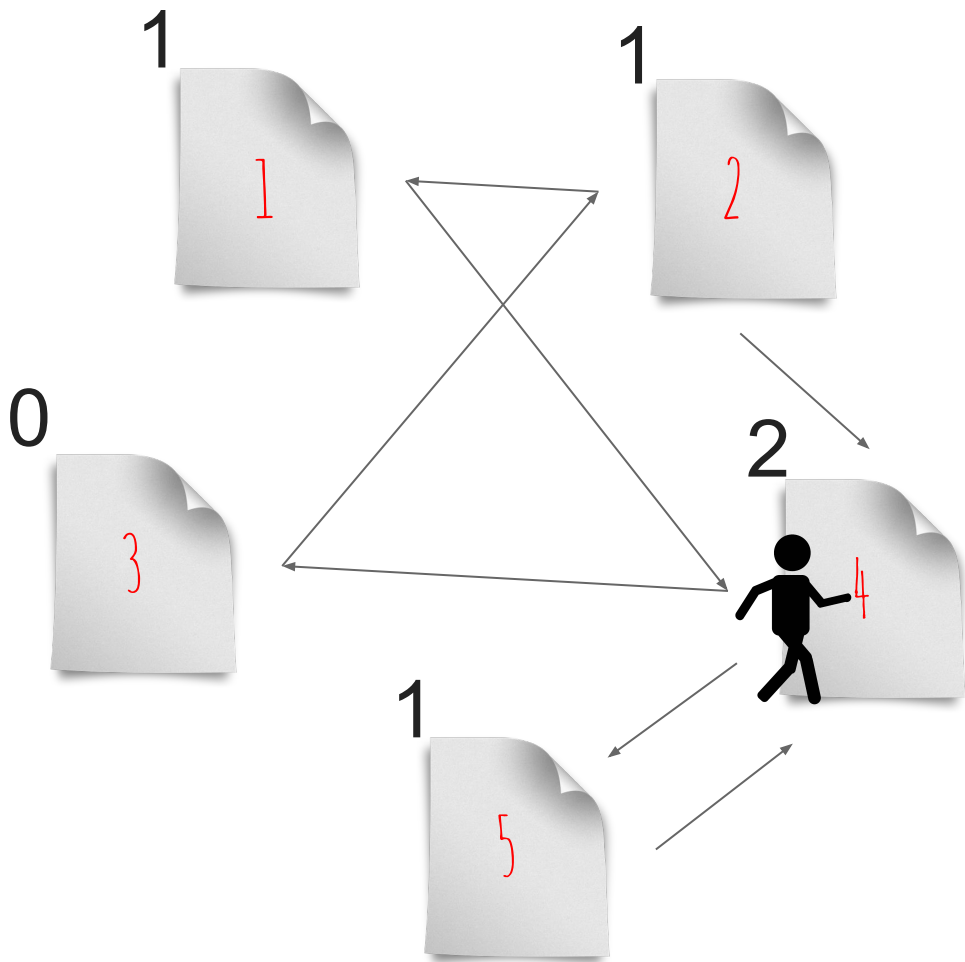


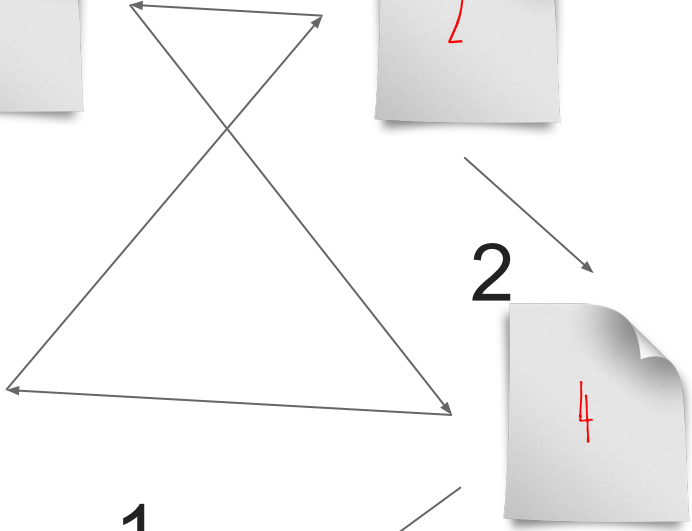
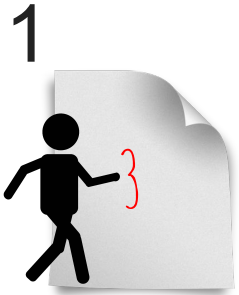


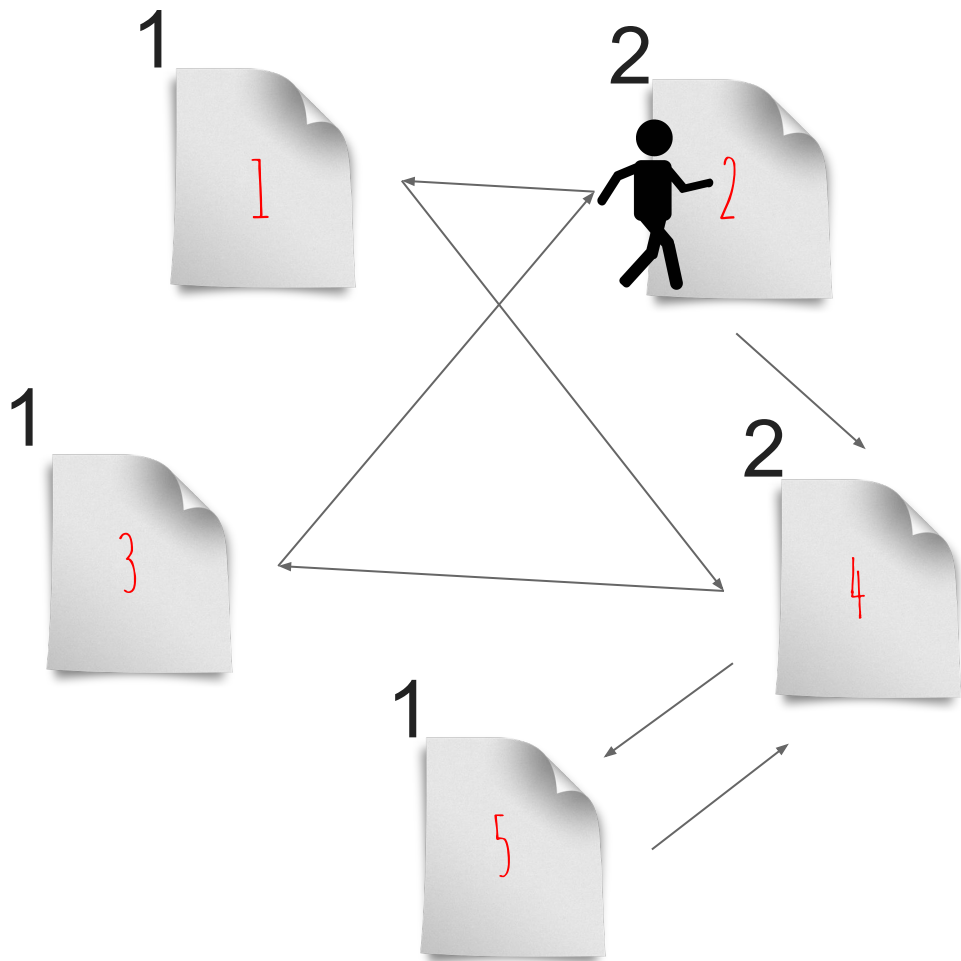


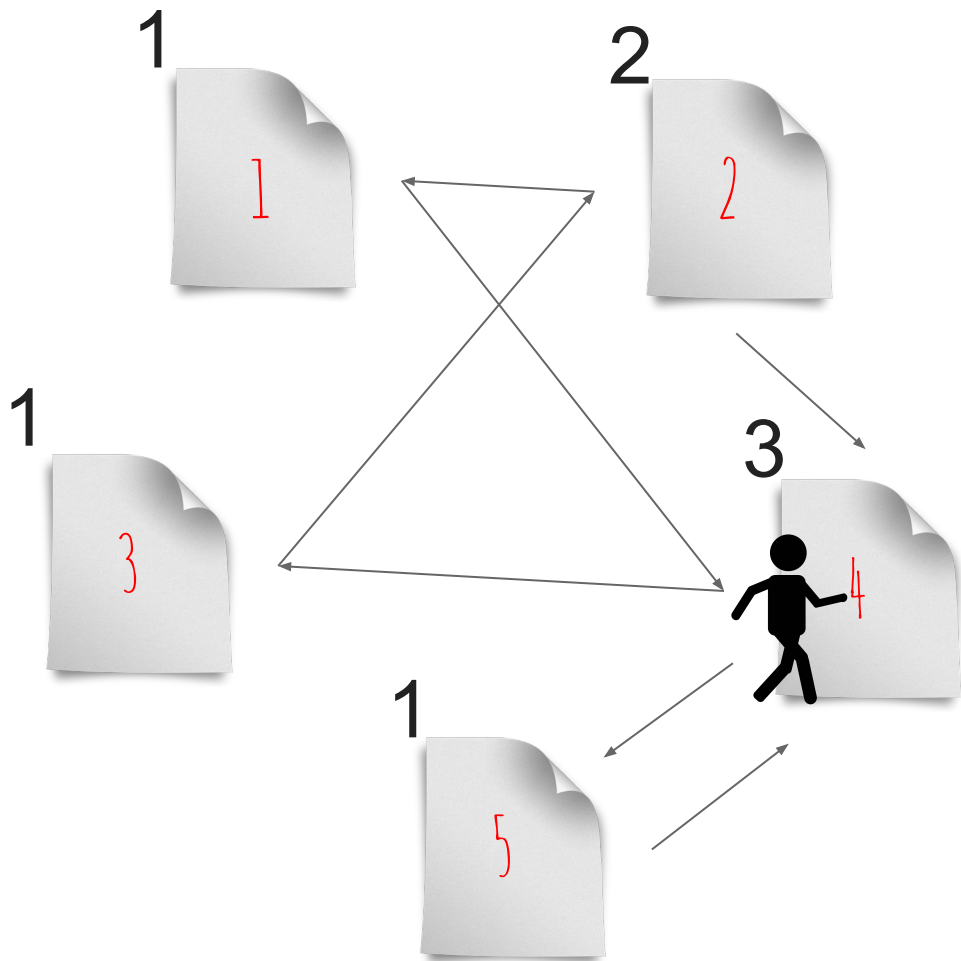


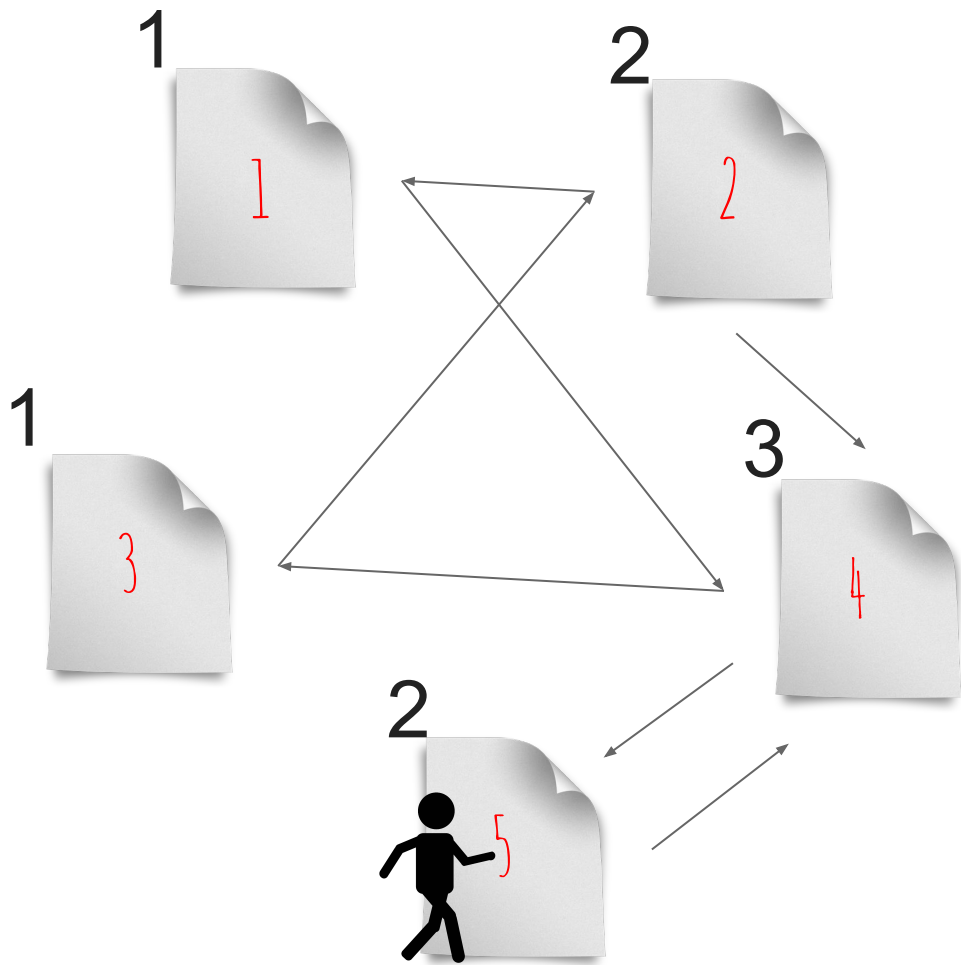




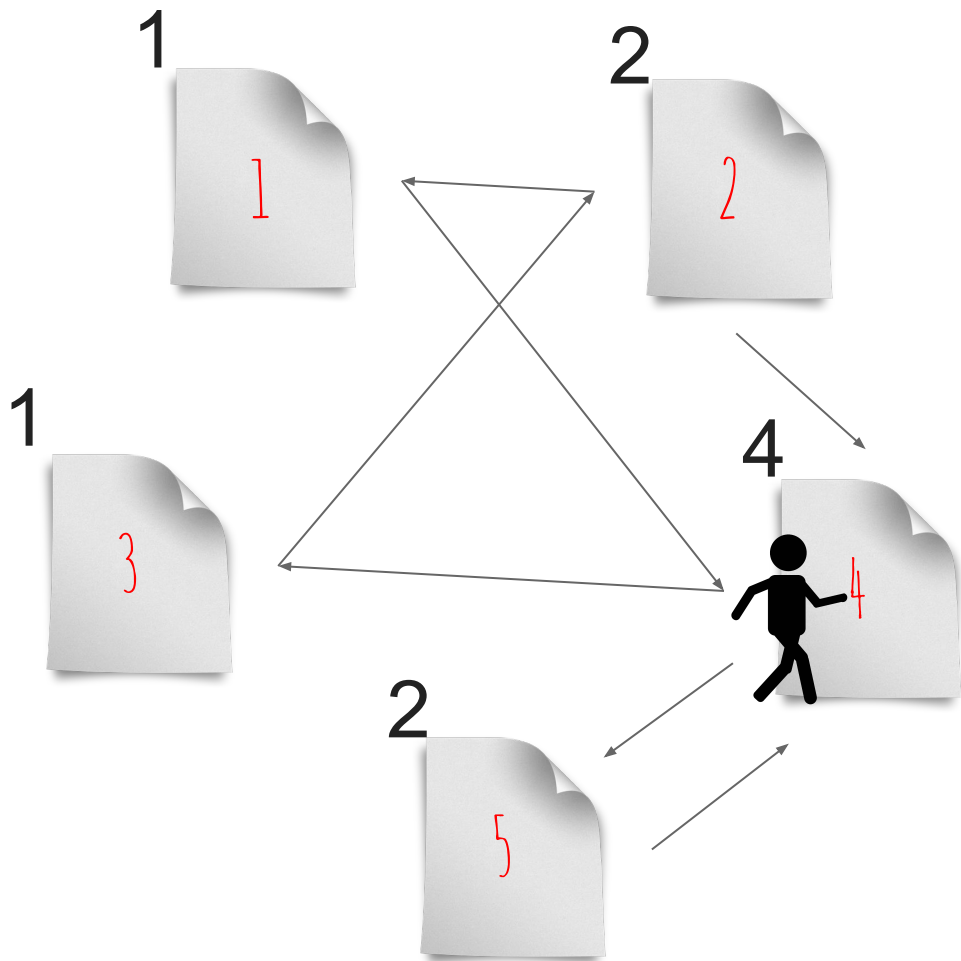


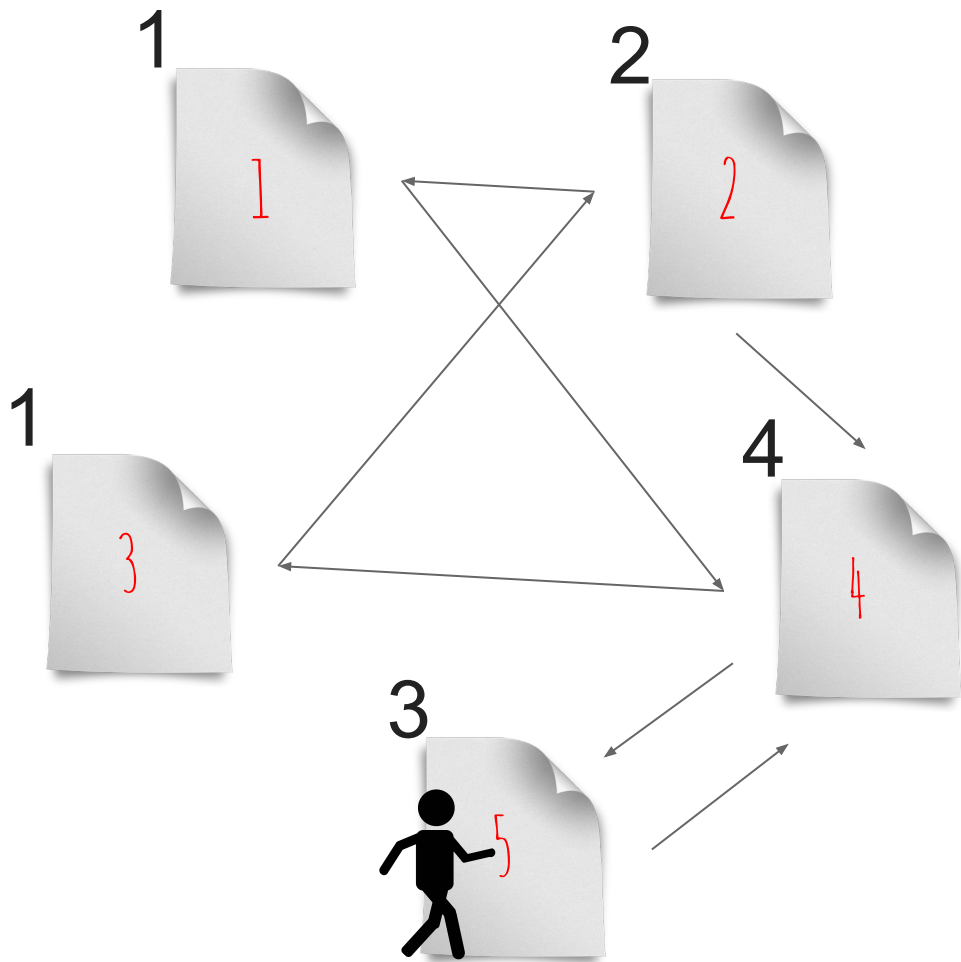


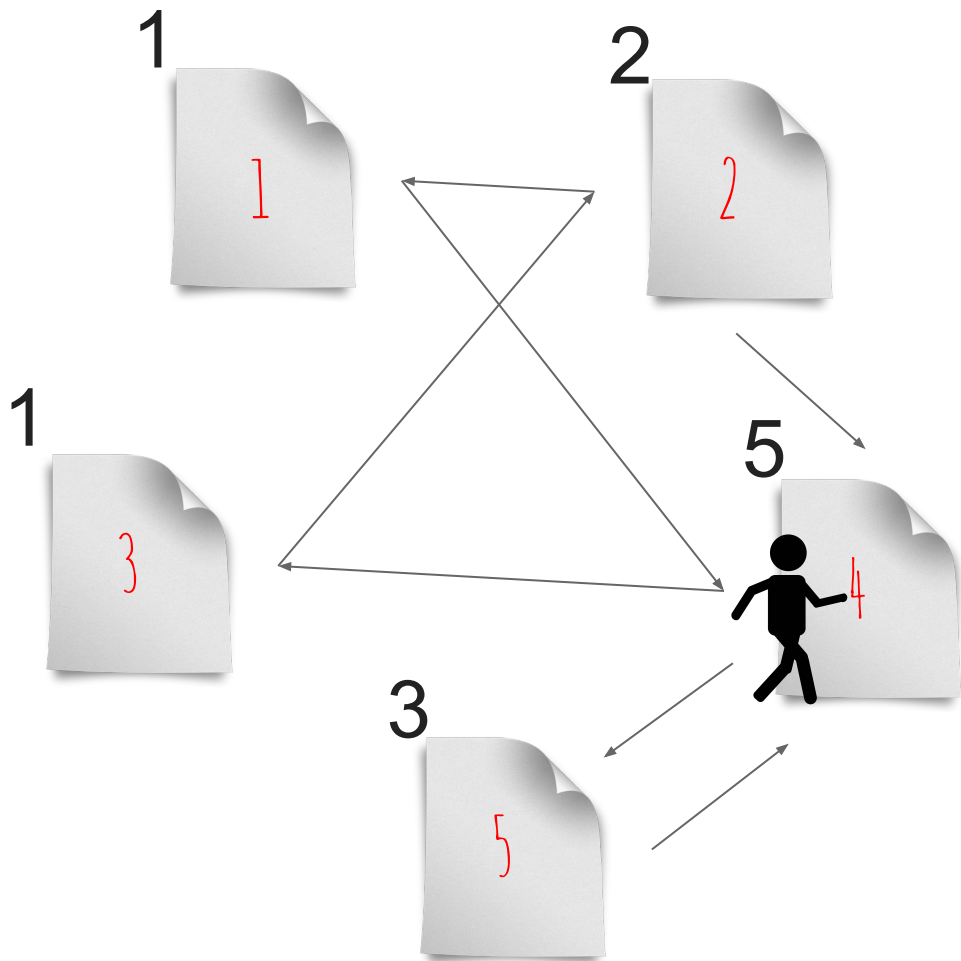


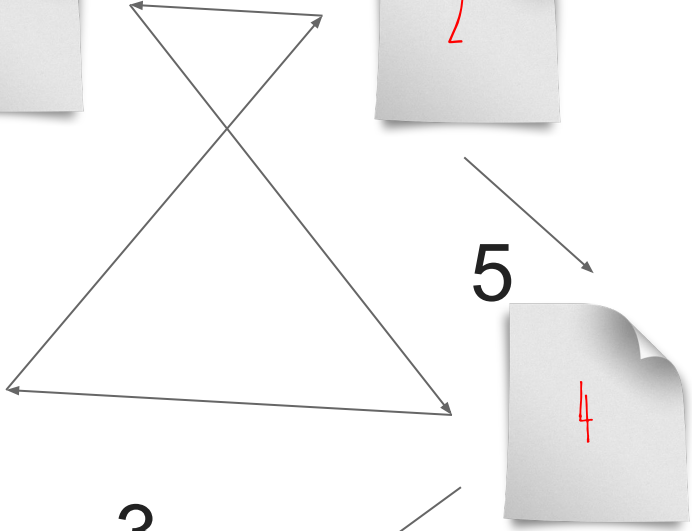
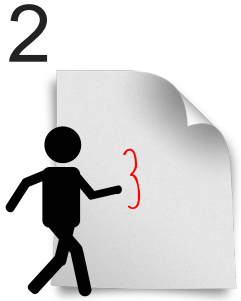


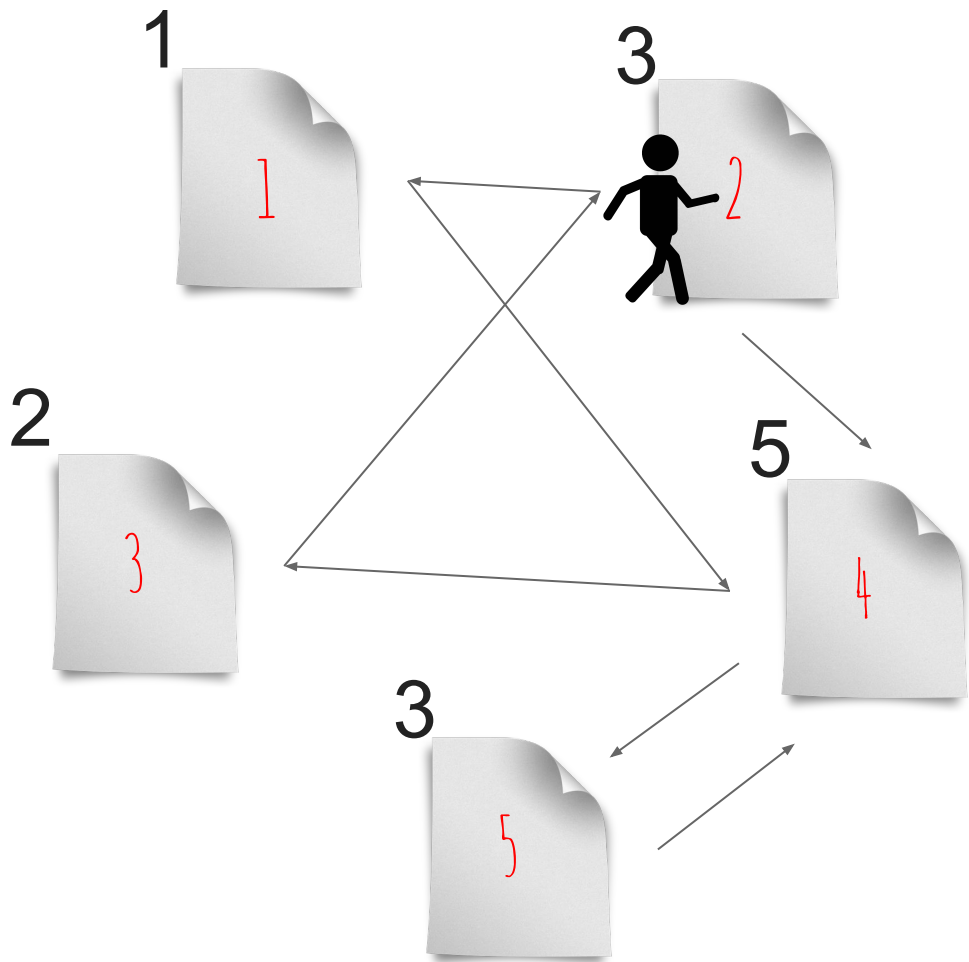


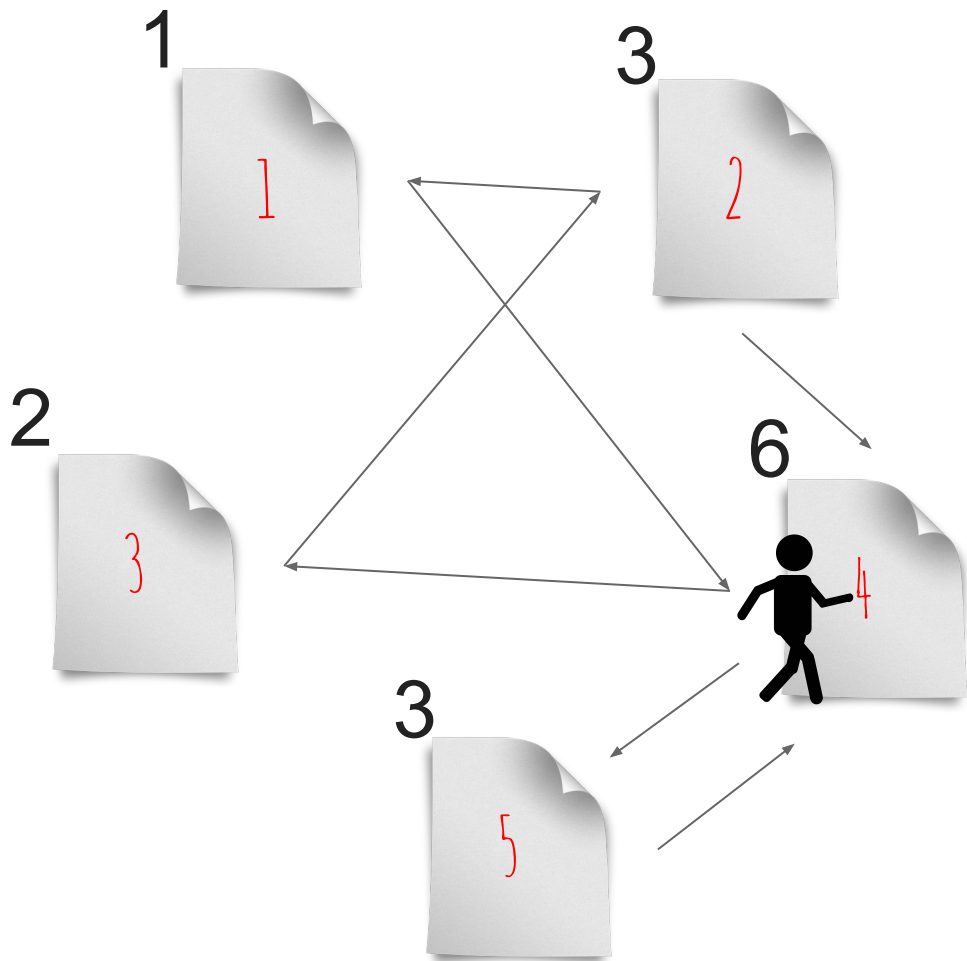


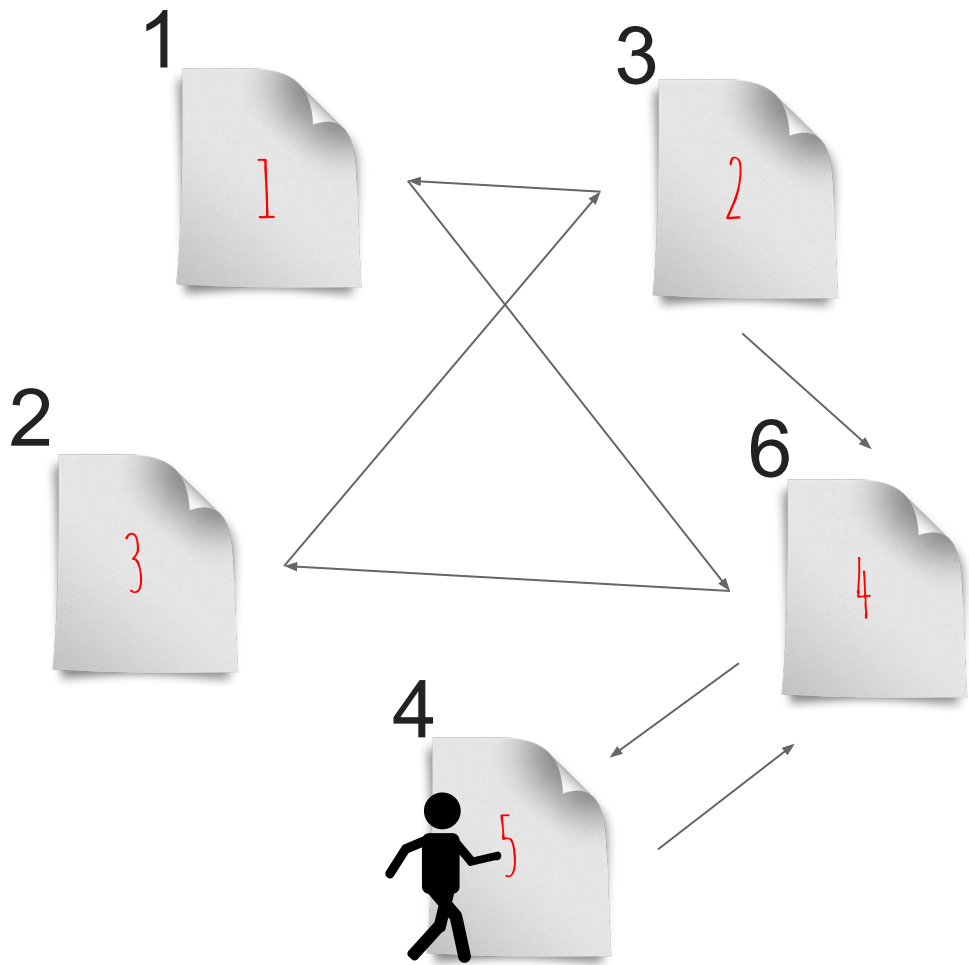


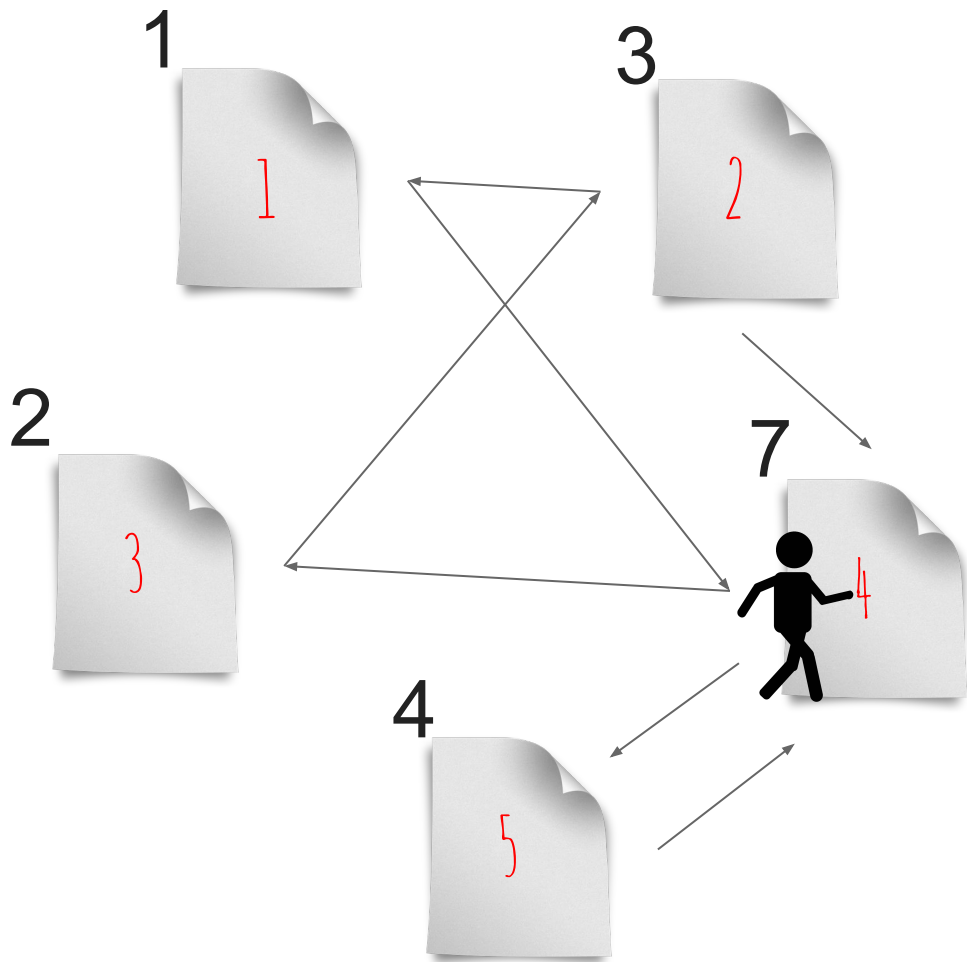




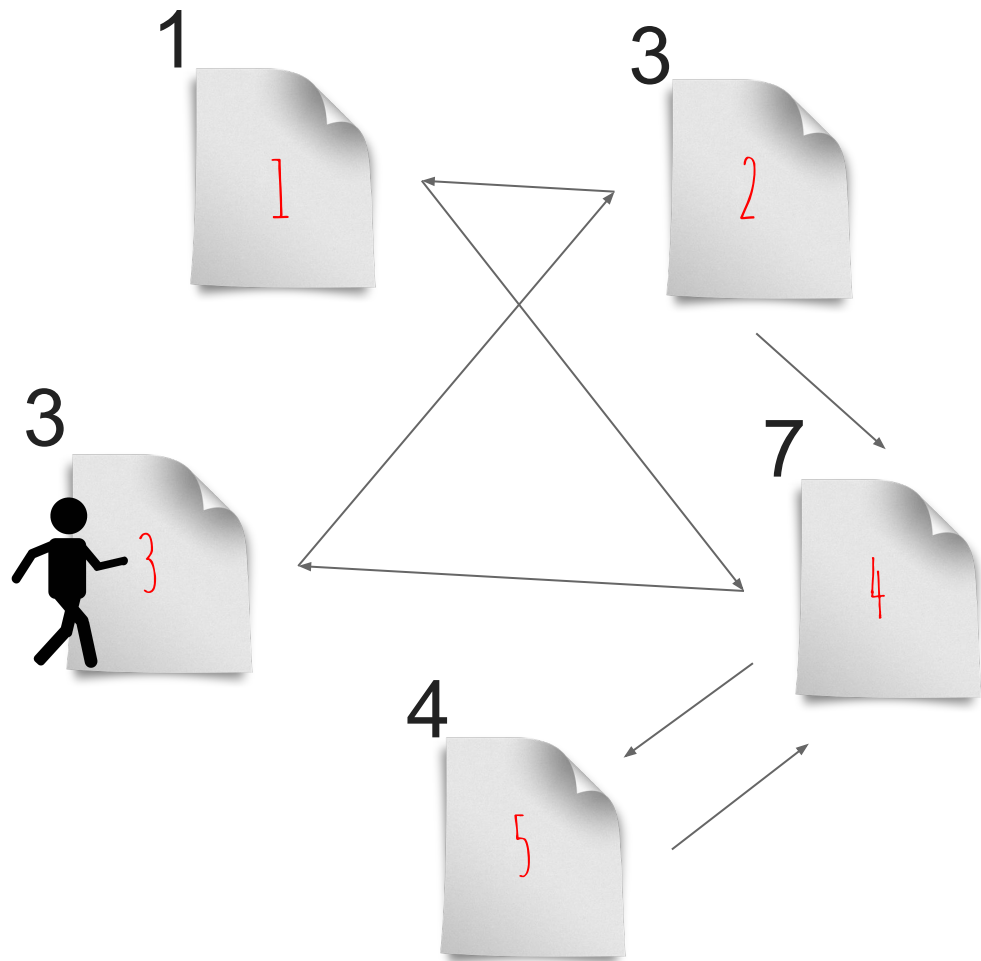




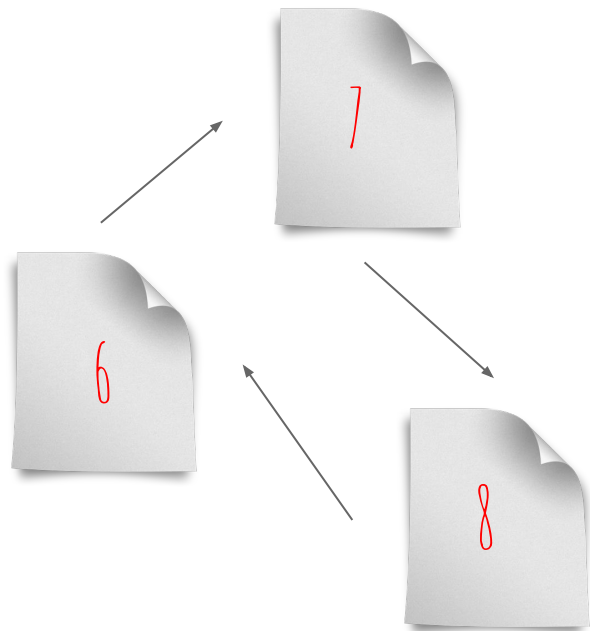
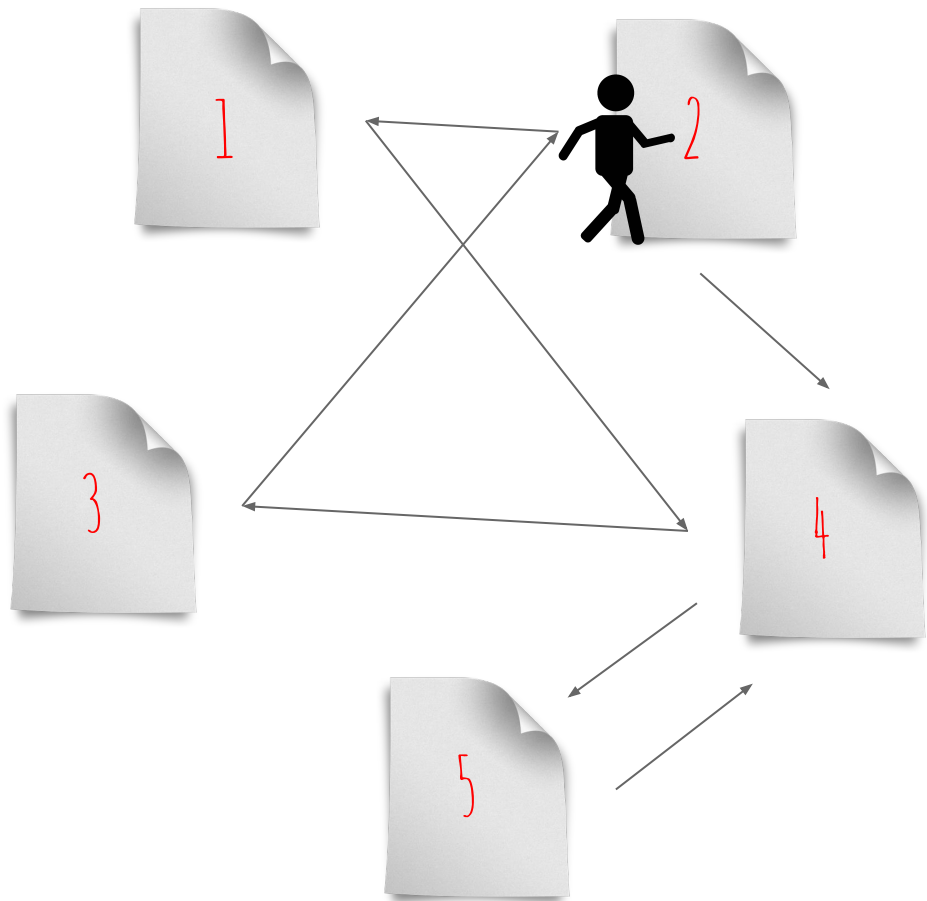


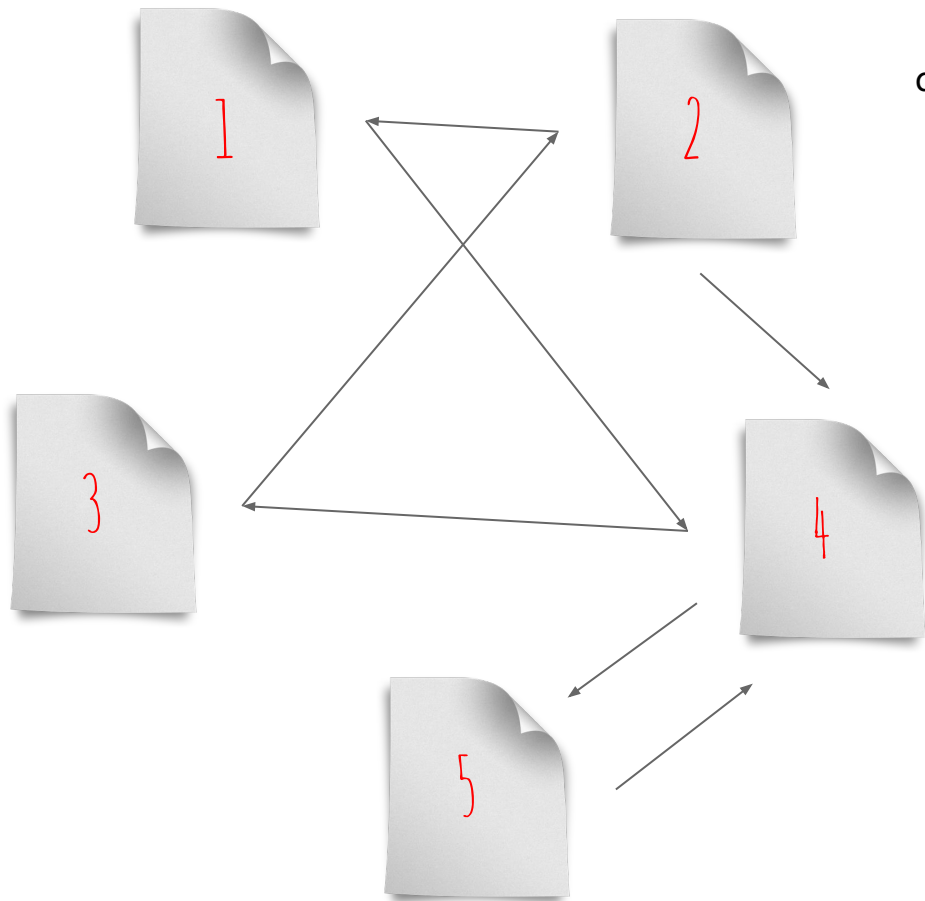




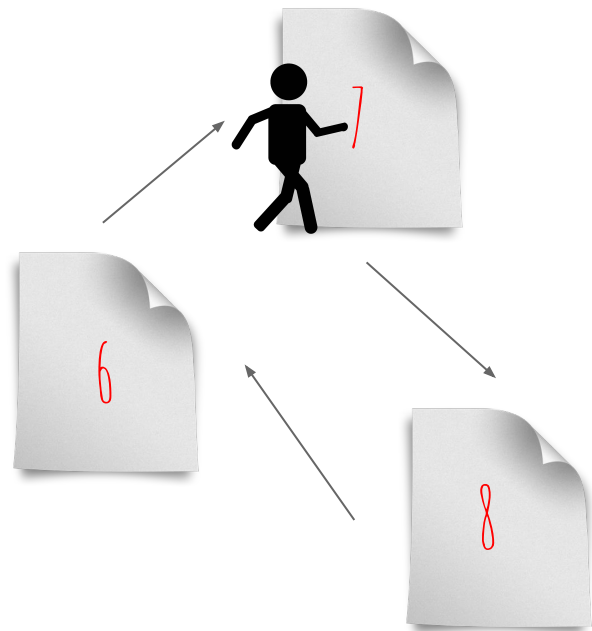


IS THE WEB CONNECTED?





$d = 0.85$



# WHAT IS BEING SAVED?

- PageRank
- Number of links per page
- Word occurrences
- Word positions
- Capitalization

# DOCID

docID				
1	cat	home	ball	...
2	ball	park	home	...
3	home	paint	people	...

# WORDID

wordID				
cat	doc1			...
ball	doc2	doc1		...
home	doc3	doc1	doc2	...

FINALLY,  
SEARCH QUERIES



- Hits
- Position
- Capitalization
- PageRank
- Anchor text

THANKS!