Au 1

# Chapter 9
# Searching The Web

# Web Search Using IR



Handwritten annotations: Aug 3:35; agent walker; ทำการ index เก็บไว้ (above Document corpus); intersec ข้อมูล Doc กับ query (mากที่แบบ แก้ไปใหม่)

Web

Spider

Document corpus

Query String

IR System

1. Page1
2. Page2
3. Page3
    .
    .

Ranked Documents

4:34

# Standard Web Search Engine Architecture

crawl the web

Check for duplicates, store the documents

DocIds

user query

create an inverted index

Show results To user

Search engine servers

Inverted index

# Challenges

ความท้าทาย ที่มีอยู่ในการ สืบทอบเว็บ

1. Distributed Data  ข้อมูลกระจายหลาย node ( หน้าที่ของ search ดึงหน้าให้ user ครบถ้วน )

2. High percentage of volatile data  เปรียนตลอดเวลา ( หน้าที่ ... ต้องตอบสนอง user ให้ทัน )

3. Large volume

4. Unstructure and redundant data
   บางอันมีโครงสร้าง / บางอันไม่มีโครงสร้าง          มันซ้ำซ้อน

5. Heterogeneous data
   เกินด้วยภาษาต่างกัน
          different languages

# Search Engines

ทำ index ไว้ตรงกลาง

## 1. Centralized Architecture
## 2. Distributed Architecture กระจายออกไป

# Centralized Architecture
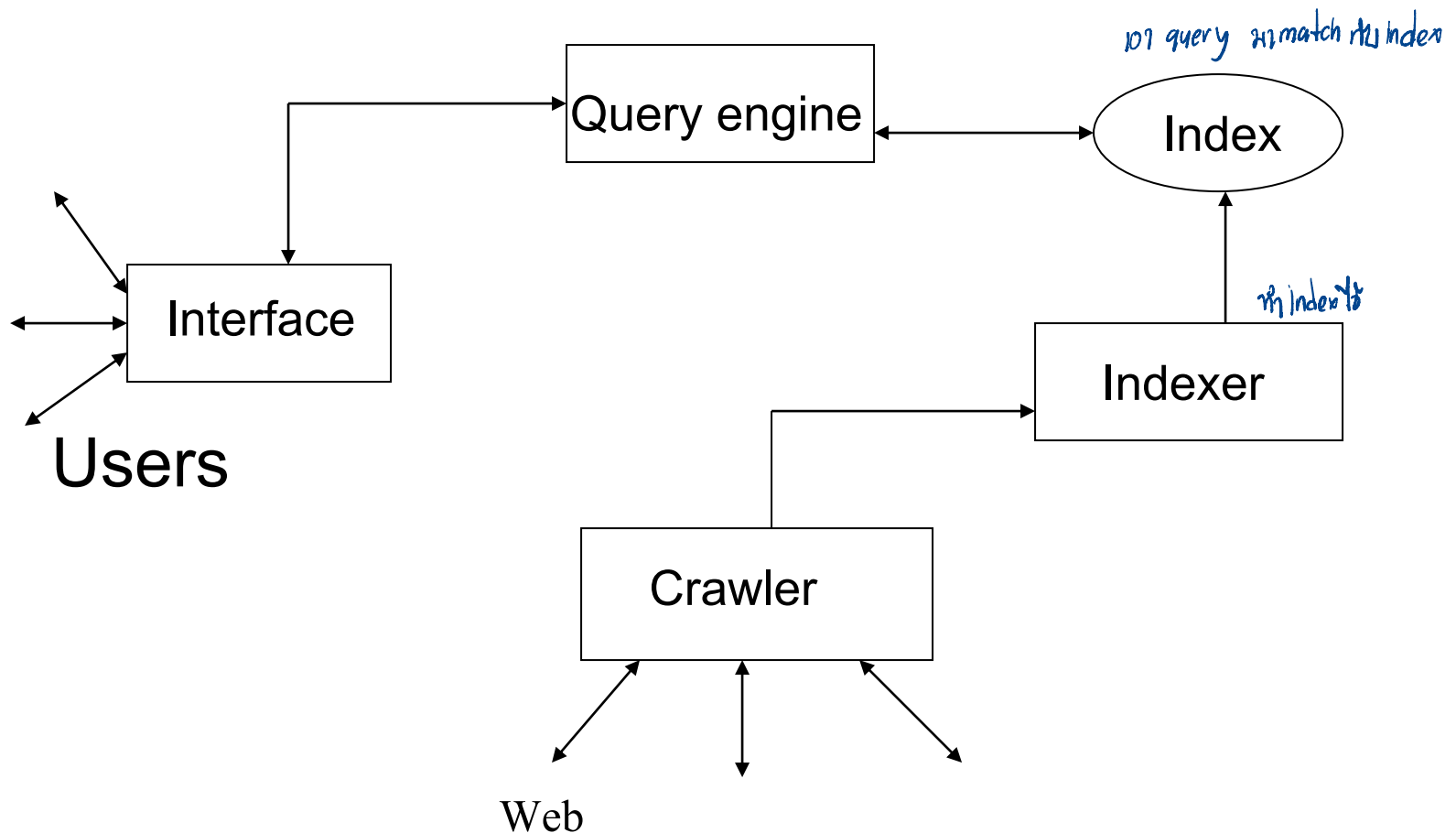# ( Crawler-indexer )

## **Definition**

  1.  Crawlers are program (software agents) that traverse the Web sending new or updated pages to a main server where they are indexed.

  2. Run on local server and send request to remote servers
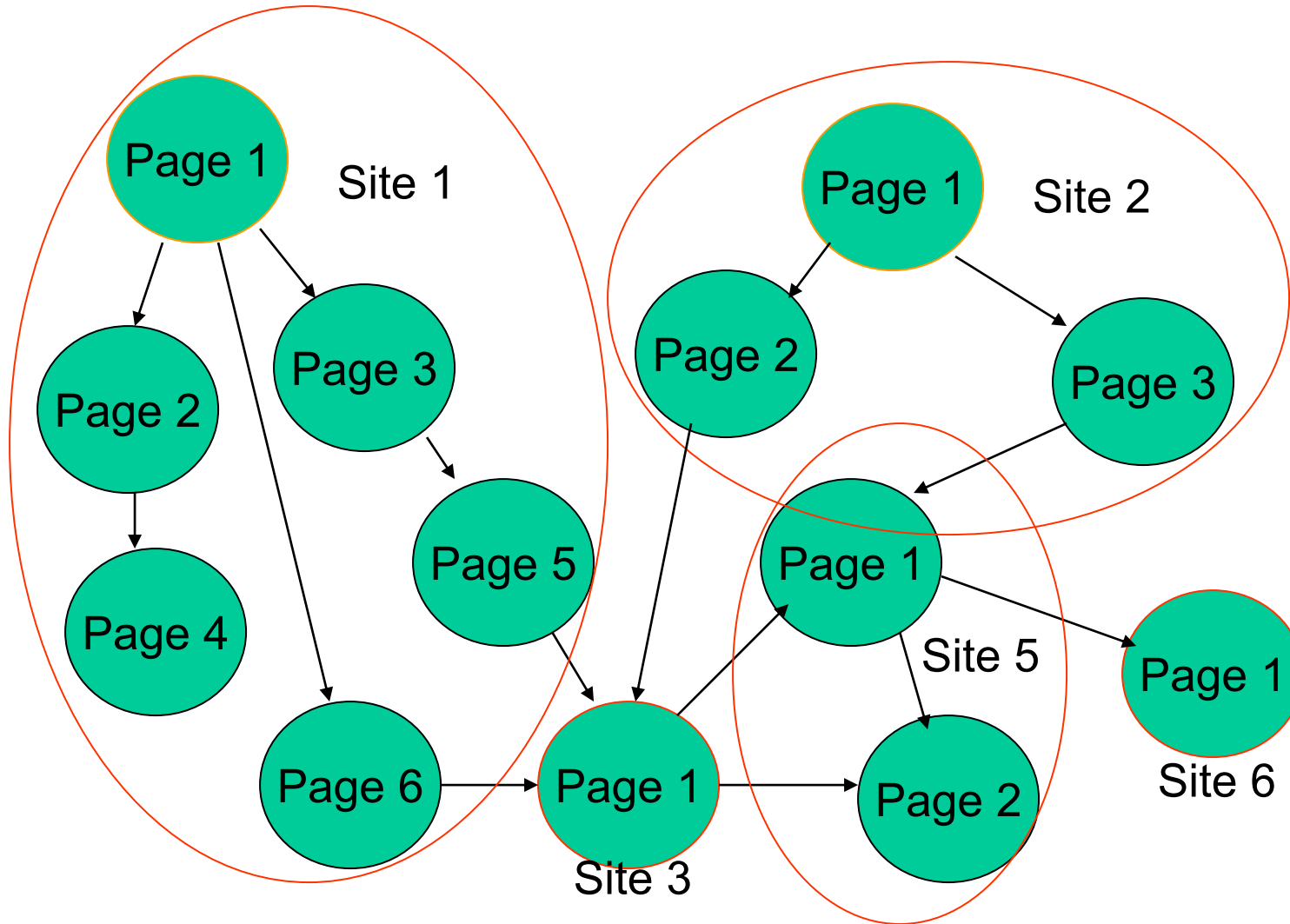
  3. Centralised use of index to answer queries

## **Name**

  Robots, Spiders, Wanderers, Walkers , Knowbot

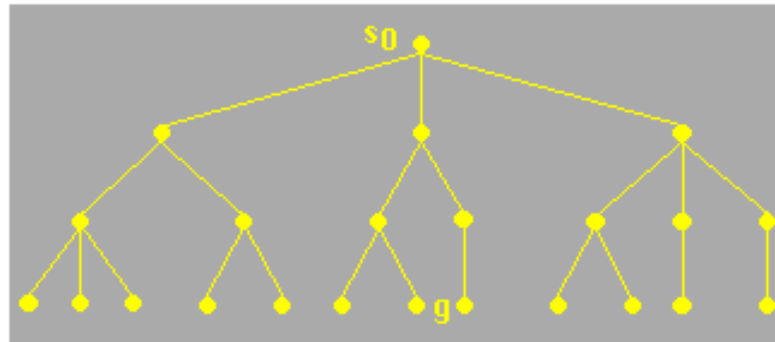# Centralized Architecture
# ( Crawler-indexer )

Query engine

Index

לפי query שמחזיר match לפי Index

Interface

Indexer

את Index מעדכן

Users

Crawler

Web

# Depth-First Crawling
## (more complex – graphs & sites)



| Site | Page |
|---:|---:|
| 1 | 1 |
| 1 | 2 |
| 1 | 4 |
| 1 | 6 |
| 1 | 3 |
| 1 | 5 |
| 3 | 1 |
| 5 | 1 |
| 6 | 1 |
| 5 | 2 |
| 2 | 1 |
| 2 | 2 |
| 2 | 3 |

8

# Depth-first search

ต้องเปรียบ tree ใช้ มากๆ

# Breadth First Crawling
# (more complex – graphs & sites)



| Site | Page |
|------|------|
| 1 | 1 |
| 2 | 1 |
| 1 | 2 |
| 1 | 6 |
| 1 | 3 |
| 2 | 2 |
| 2 | 3 |
| 1 | 4 |
| 3 | 1 |
| 1 | 5 |
| 5 | 1 |
| 5 | 2 |
| 6 | 1 |

# Breadth-first search

# Centralized Architecture ( Crawler-indexer )

*24.58*

*ดี : ง่าย*

## Problem

1. Volumn of the data

*Traffic ∞ data*

2. Traffic (Crawler retrieve entire object) *ทั้งหมด ออก-มา*

3. High load at Web Servers

↓

*มีเทคนิค ในการแบ่ง ลื้อ Die*

# Distributed Architecture ( Harvest )

## Definition

*data ทิน*

*แหล่งความรู้*

*1 Broker : หลาย/1 Gather*

   1. <u>Gatherers</u> collects and extracts indexing information form one or more Web servers at periodic time

*→ ส่ง index จากการหาง broker ในแต่ละ ความรู้ได้*

*query ที่หา = ตอบปาเอง*

*↪ ไม่ตอบ = โยนไปนักการอื่นที่ตอบ (Broker อื่น)*

   2. <u>Brokers</u>  *หาง*

*หลายนักการต่อ1 ตัวสามน่อง ความรู้ เป็น ของตัวเอง  มีหลาย*
*ชนิดได้*

*ถ.มี*

     -Provide <u>indexing</u> mechanism and query interface to data gathered

*index*

     -Retrieve <u>information</u> from gatherers or other brokers, updating incrementally their indices

# Distributed Architecture
# ( Harvest architecture )



Replication Manager

Broker

Broker

User

Gatherer

Object cache

Web site