# LC029 정보검색

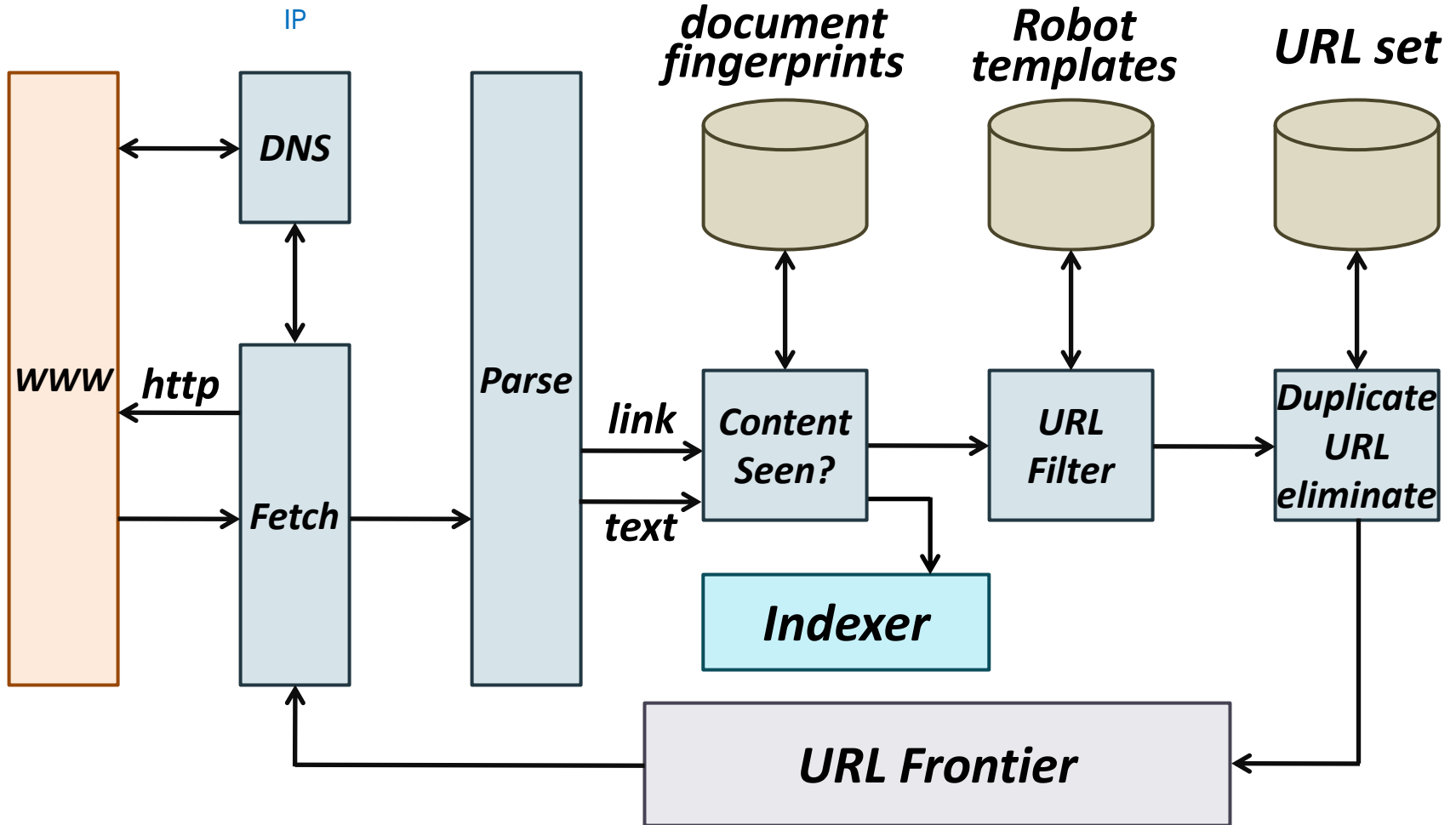## Chapter 20 : Web crawling and indexes

1. A Crawler
2. Features of a Crawler
3. Distributing the crawler
4. URL Frontier
5. Connectivity Servers
6. Distributing Indexes

# A Crawler

# Basic Crawler Architecture



WWW

DNS

IP

http

Fetch

Parse

link

text

document fingerprints

Content Seen?

Indexer

Robot templates

URL Filter

URL set

Duplicate URL eliminate

URL Frontier

# Basic crawler operation

- Pick a URL <u>from the frontier</u>.    Begin with a seed set.
- Fetch the document at that URL.
- Parse the fetched document.
    - Extract links to other docs (URLs)
- Check if the document has content already seen.
    - If not, add to indexes.
- For each extracted URL    e.g. only crawl .edu obey robots.txt, etc
    - Ensure it passes certain URL <u>filter tests.</u>
    - Check if it is already in the frontier (duplicate URL elimination).

# DNS (Domain Name Server)

- Provide a lookup service on the internet.
  - Given a URL, retrieve its IP address. URL     IP
- Service provided by a distributed set of servers. Therefore, lookup latencies can be high (even seconds). DNS
  - DNS caching is used.
- Common DNS implementations are *blocking*. (block)
  - Once a request is made to DNS, other requests are blocked until the first request is completed.
  - Most web crawlers implement their own DNS resolver as a component of the crawler.

# Parsing: URL normalization

- When a fetched document is parsed, some of the extracted links are *relative* URLs.

- e.g. at http://en.wikipedia.org/wiki/main.htm

> **<A href="other/page.htm">Other Page</A>**

- Extracted link is "other/page.htm" which is the same as the absolute URL
  http://en.wikipedia.org/wiki/other/page.htm.

- During parsing, we must normalize such relative URLs.

# Content seen?

- Duplication is widespread on the web.

- If the page just fetched is already in the index, do not further process it.

- This is verified using document fingerprints and shingles.

# Filters and robots.txt

- **Filters** URL
  - Determine URLs to be excluded from the frontier.
  - Regular expressions are used for URLs to be crawled/not.
- **robots.txt**
  - **Robot Exclusion Protocol** for giving spiders ("robots") limited access to a website, originally from 1994.
    For more, see [www.robotstxt.org/robotstxt.html](http://www.robotstxt.org/robotstxt.html).
  - Once a robots.txt file is fetched from a site, do not fetch it repeatedly. (doing so burns bandwidth, hits web server)
  - Therefore, we do cache robots.txt files.

# robots.txt : example

- For a URL, create a file **URL/robots.txt**.

```
# robots.txt for http://www.example.com/

User-agent: *
Disallow: /cgi-bin/

User-agent: searchengine
Disallow:
```

- No robot should visit any pages in **URL/cgi-bin/**, except a web robot called "searchengine".

# Housekeeping Tasks

- Typically performed by a dedicated thread.

- It wakes up once every few seconds to log crawl progress statistics (URLs crawled, frontier size, etc)

- Periodically, it takes a snapshot of the crawler's state (say, URL frontier) and stores it on a disk.

- On a catastrophic failure, the crawling is restarted from the most recent snapshot (checkpoint).

# Features of a Crawler

# Features a crawler should provide

- Be **<u>Polite</u>**
Respect explicit and implicit politeness considerations.
  - **Explicit politeness**

    Specifications from webmasters on what portions of site can be crawled.

    - For this specification, **robots.txt** is used.
  - **Implicit politeness**

    Even with no specification, avoid hitting any site too often.                   DDOS

    - Only one connection should be open to any given host at a time.
    - A waiting time of a few seconds should occur between successive requests to a host.

# Features a crawler should provide

- Be **Robust**
  Be immune to **spider traps** and other malicious behavior from web servers. :

- Be **Scalable**
  Designed to increase the crawl rate by adding more machines.

- Be capable of **Distributed operation**
  Designed to run on multiple distributed machines.

- **Performance/Efficiency**
  Permit full and efficient use of available processing and network resources.

# Features a crawler should provide

- Fetch pages of **<u>higher quality first</u>**.      fetch

- **<u>Continuous operation</u>**
  Continue fetching fresh copies of a previously fetched page, so that the search engine's index is fairly current.

- **<u>Extensible</u>**
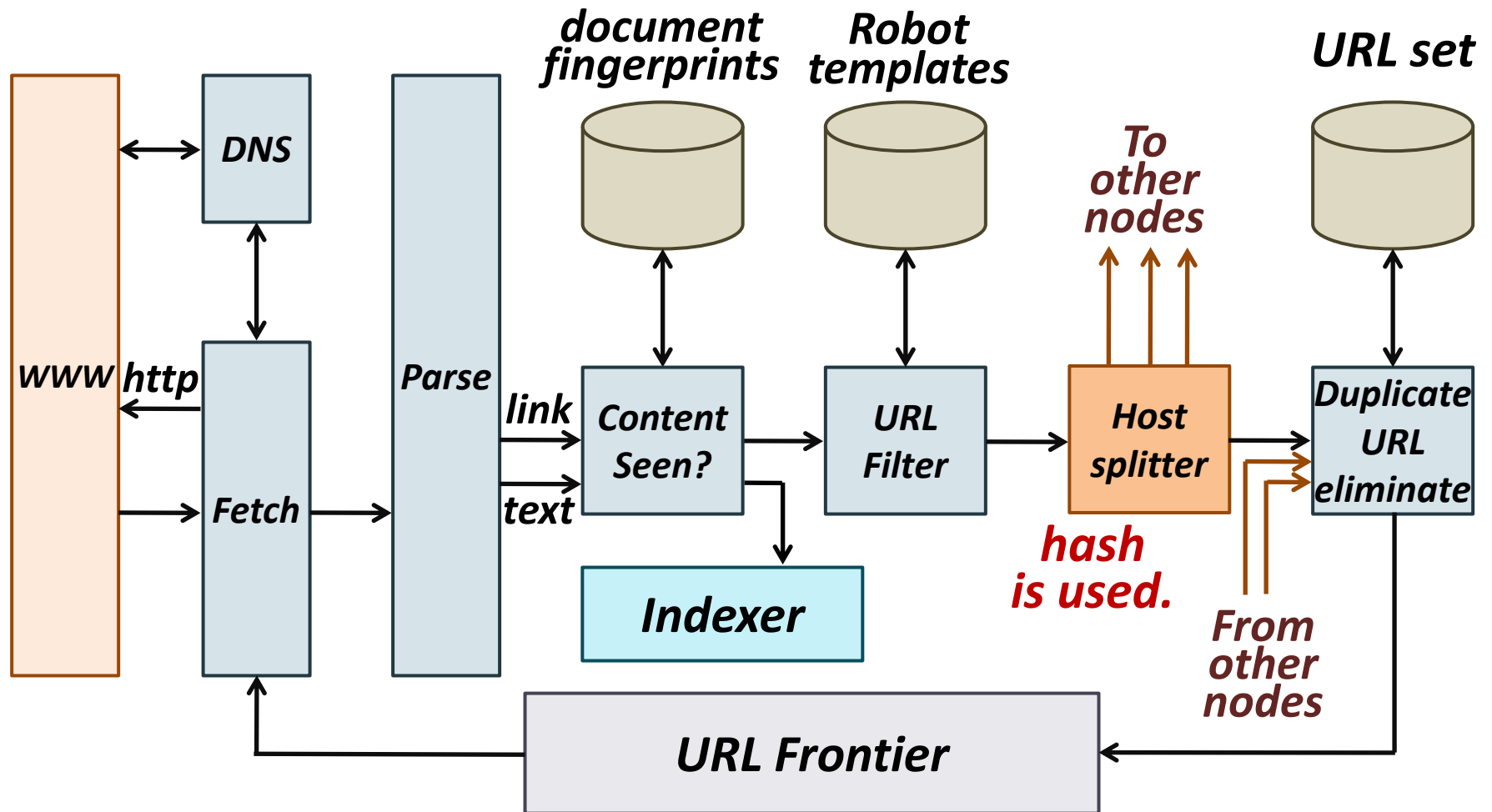  Adapt to new data formats, protocols, etc.

# Distributing the crawler

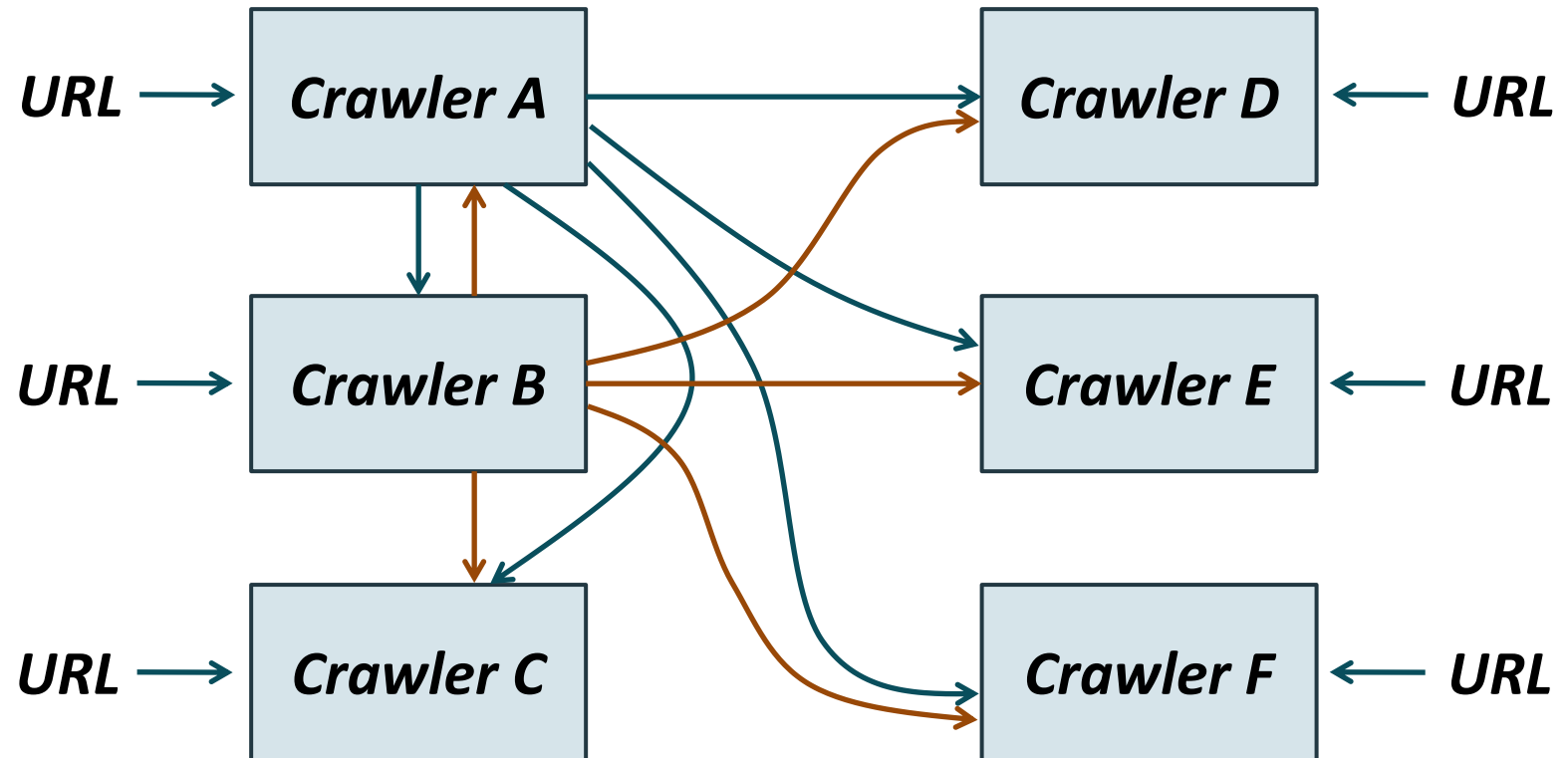# Distributing the crawler

- Run multiple crawl threads, under different processes, potentially at different nodes.

    - Geographically distributed nodes

- **Host Splitter** does partition hosts being crawled into nodes. <span>URL</span>

    - Hash used for partition.

- Distributing the crawler is essential for **scaling**.

# Distributing the crawler



document fingerprints

Robot templates

URL set

DNS

To other nodes

www http

Parse

link

Content Seen?

URL Filter

Host splitter

Duplicate URL eliminate

Fetch

text

hash is used.

From other nodes

Indexer

URL Frontier

# Distributing the crawler

# URL Frontier

# URL Frontier

- Contain URLs to be fetched in the current crawl.

    - Can include multiple pages from the same host.

- Provide a URL in some order when a crawler thread seeks one.

    - Must avoid trying to fetch them all at the same time.

    - Must try to keep all crawling threads busy.

# URL Frontier: two main considerations

- **Freshness**

  Crawl high-quality pages more often than others.
  - Pages (such as news sites) whose content changes often get high priority.

- **Politeness**

  Do not hit a web server too frequently.

  - Even if we restrict only one thread to fetch from a host, the crawler can hit it repeatedly.

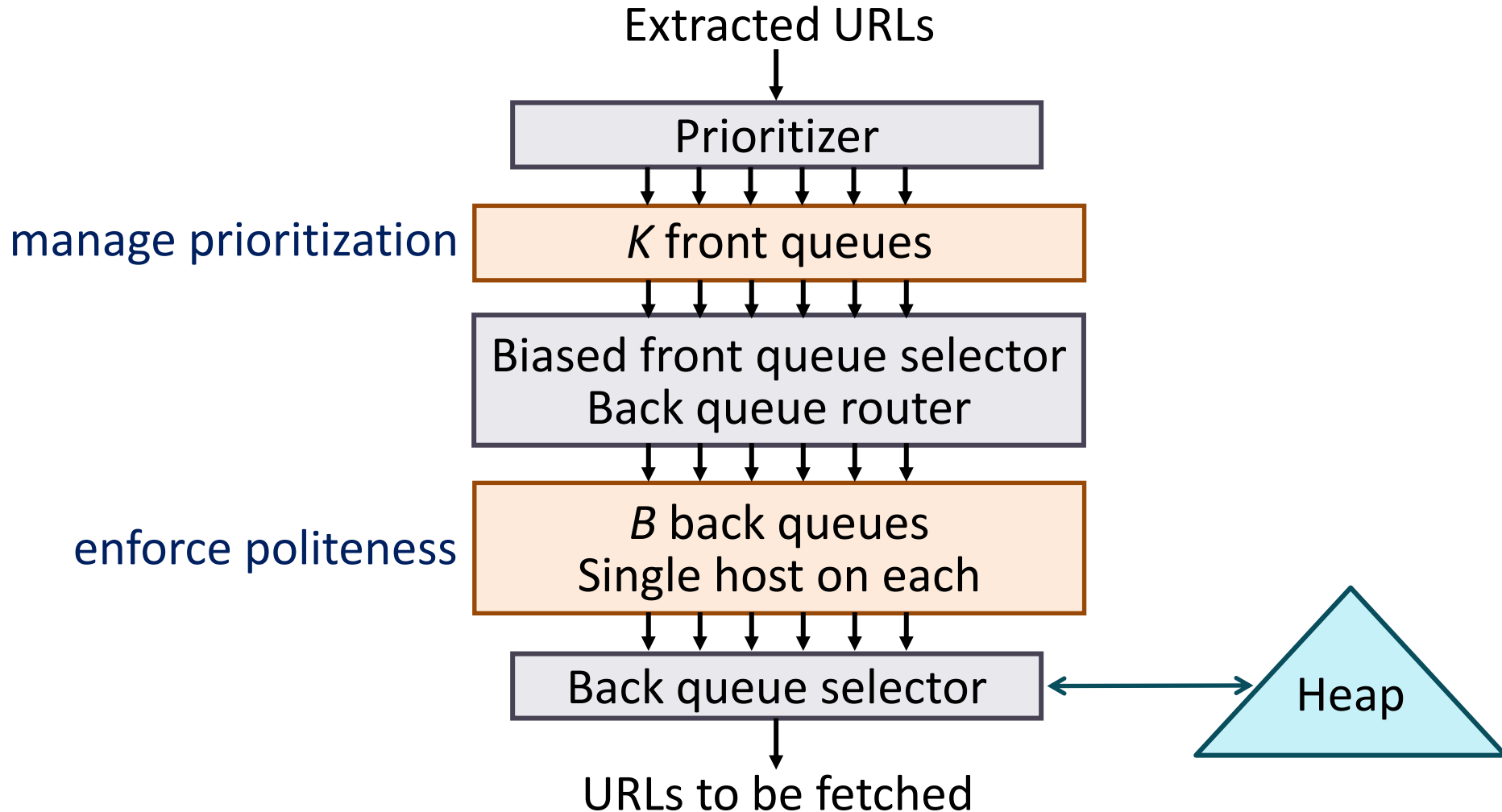  - So, a waiting time of a few seconds is inserted between successive requests to a host.

# Implementation of URL Frontier

- Simple priority queue fails because these goals may conflict each other.

    - Crawl high-quality pages more often than others.
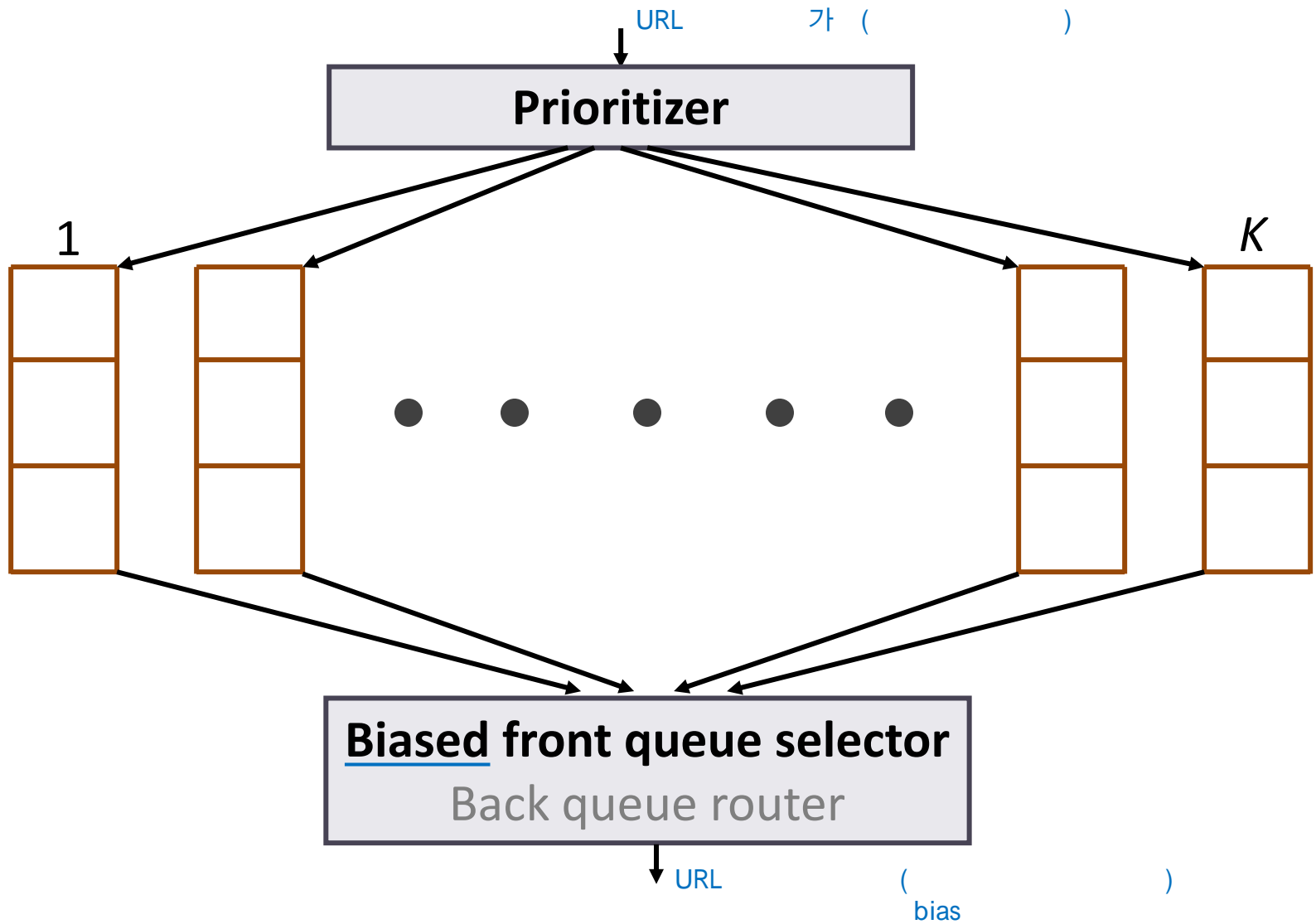    - Do not hit a web server too frequently.

    Note that many links out of a page go to its own site, creating a burst of accesses to that site.

# Implementation of URL Frontier

Extracted URLs

Prioritizer

manage prioritization — $K$ front queues

Biased front queue selector
Back queue router

enforce politeness — $B$ back queues
Single host on each

Back queue selector ⟷ Heap

URLs to be fetched

# Front Queues

# Front Queues

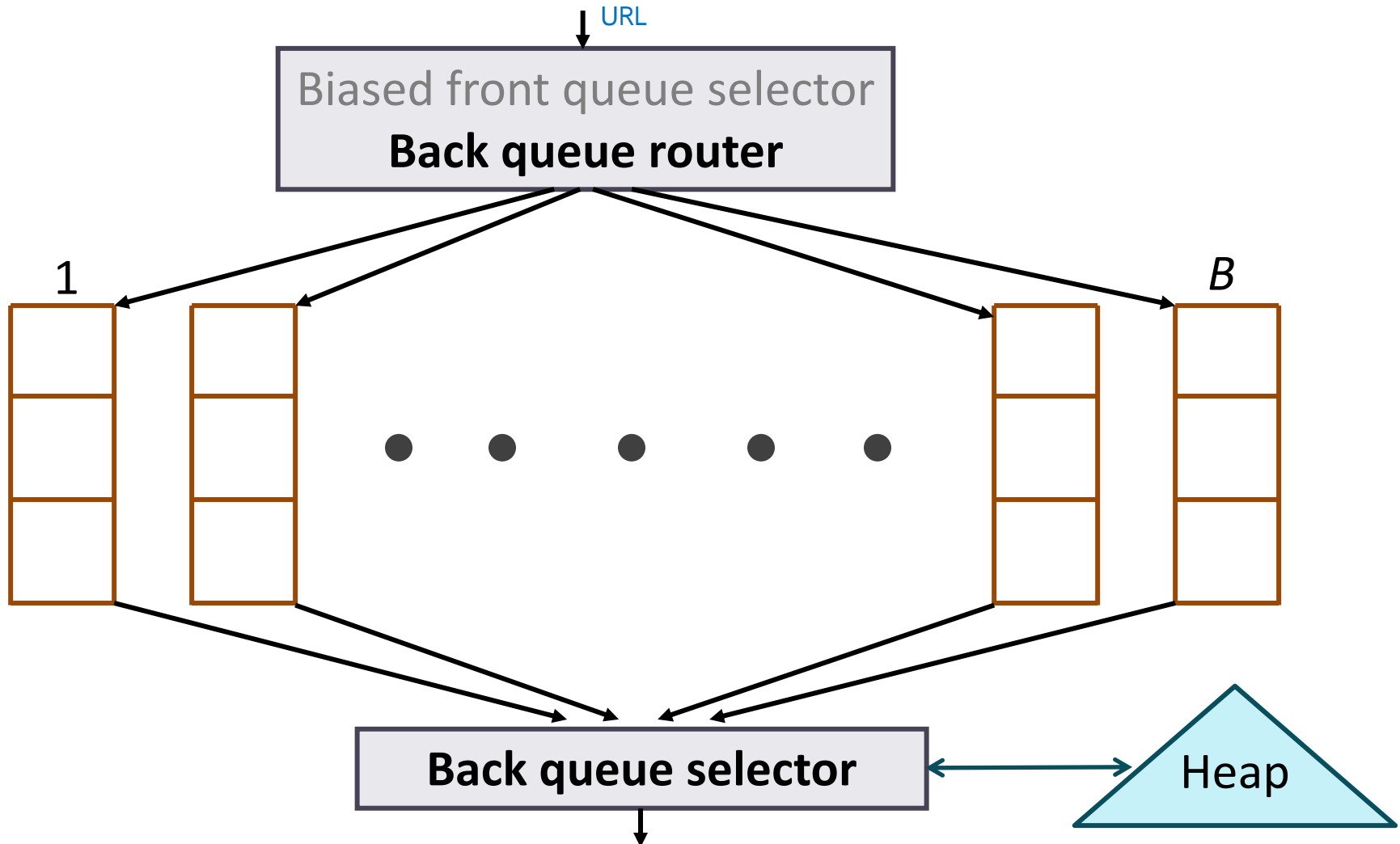- Prioritizer assigns to URL an integer priority between 1 and *K.*

    - Appends a new URL to the corresponding queue.

- Heuristics for assigning priority

    - Refresh rate sampled from previous crawls.

    - Application-specific
      e.g. crawl news sites more often.

# Biased Front Queue Selector

- When a <span style="color:green">back queue</span> requests a URL,
  picks a <span style="color:blue">front queue</span> from which to pull a URL.

- This choice can be round robin biased to queues of higher priority, or some more sophisticated variant.

  - Can be randomized

# Back Queues

# Back Queues

- Each back queue is kept non-empty while the crawl is in progress.

- Each back queue only contains URLs from a single host.

  - Maintain a table from hosts to back queues.

| Host name | Back queue |
|---|---|
| stanford.edu | 23 |
| microsoft.com | 47 |
| sungshin.ac.kr | *12* |

# The Heap

- One entry for each back queue.

- The entry is the earliest time $t_e$ at which the host corresponding to the back queue can be hit again.

- This earliest time is determined from

  - Last access to that host.

  - Any time buffer heuristic we choose.

:

# Back Queue Processing

- When a crawler thread seeking a URL to crawl,
  - Extracts the root of the heap.
  - Wait until the corresponding time entry $t_e$.
- Fetches URL at head of corresponding back queue $q$.
- Create a new heap entry for the URL.
- Checks if queue $q$ is now empty.
  - if so, pulls a URL $v$ from front queues.
  - if there is a back queue $q'$ for $v$'s host, append $v$ to $q'$ and pull another URL from front queues, repeat until $v$ for $q$ is found. URL
    URL URL , URL
  - add $v$ to $q$

# Number of Back Queues *B*

- Keep all threads busy while respecting politeness.

- Mercator recommendation
  Three times as many back queues as crawler threads.

# Size of URL Frontier

- On a web-scale crawl, the URL frontier may grow too big to reside in memory.

- Solution
  - Let most of the URL frontier reside on disk.
  - A portion of each queue is kept in memory, with more brought in from disk as it is drained in memory.
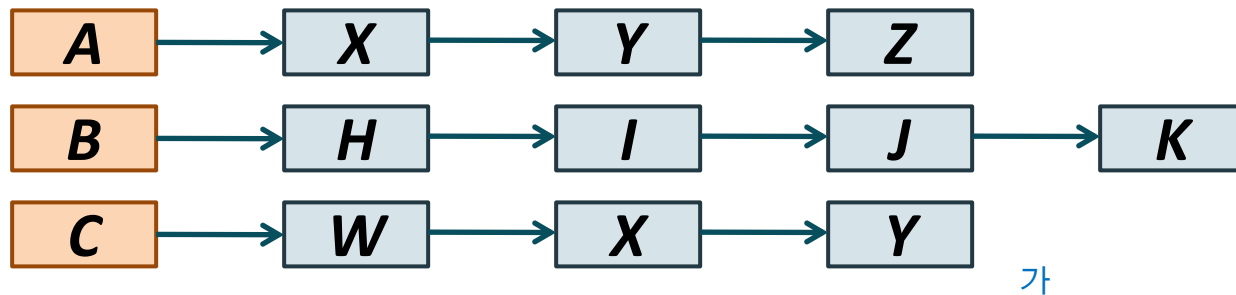  
  ( )

# Connectivity Servers

# Connectivity Server

- Support for fast queries on the web graph
    - Which URLs point to a given URL?
    - Which URLs does a given URL point to?

    *adjacency table*

    | A | → | X | → | Y | → | Z | |
    |---|---|---|---|---|---|---|---|
    | B | → | H | → | I | → | J | → K |
    | C | → | W | → | X | → | Y | |

- Applications
    - Web graph analysis for sophisticated crawl optimization
    - Link analysis

# Connectivity Server
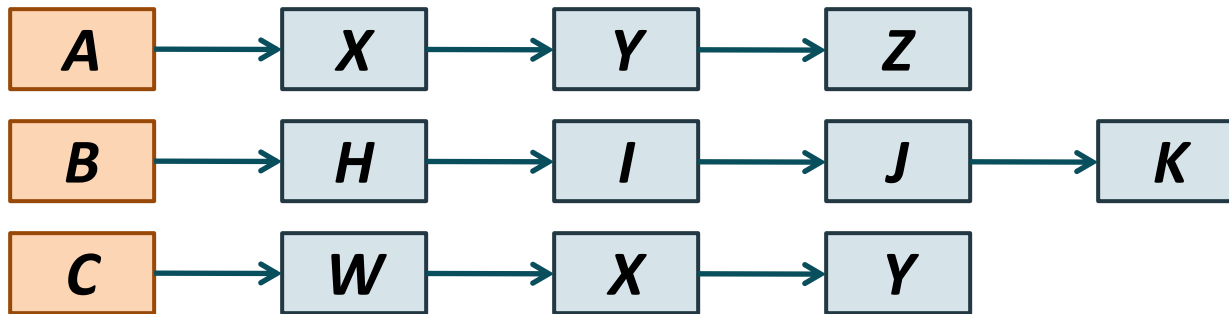
- Store the adjacency table in memory
  - From URL to outlinks
  - From URL to inlinks

- Memory requirement
  - Assume that the Web had 4 billion pages, each with ten links to other pages.
  - Assuming each URL represented by an integer, we need 32 bit integer per node.

    $4 \times 10^9 \times 10 \times 4 \times 2 = 3.2 \times 10^{11} = 3.2 \times 10^2$ TB

- Thus, compressing the adjacency table is critical.

# Compression of Adjacency Table

- Properties exploited in compression:
  - **Similarity** (between lists)

| A | → | X | → | Y | → | Z | | |
|---|---|---|---|---|---|---|---|---|
| B | → | H | → | I | → | J | → | K |
| C | → | W | → | X | → | Y | | |

From URL to outlinks

  - **Locality** (many links from a page go to "nearby" pages)
    - Use gap encodings in sorted lists

# Compression of Adjacency Table

| | |
|---|---|
| **nate.com/a.htm** | nate.com/w.htm, naver.com/x.htm |
| **nate.com/b.htm** | nate.com/cafe/boss.htm, nate.com/y.htm |
| **nate.com/c.htm** | nate.com/cafe/boss.htm, naver.com/x.htm |
| **nate.com/d.htm** | nate.com/w.htm, nate.com/y.htm, ever.com/z.htm |

Adjacency Table
(From URL to outlinks)

# Compression of Adjacency Table

- Consider lexicographically ordered list of all URLs.

| 1 | www.stanford.edu/alchemy |
|---|---|
| 2 | www.stanford.edu/biology |
| 3 | www.stanford.edu/biology/plant |
| 4 | www.stanford.edu/biology/plant/copyright |
| 5 | www.stanford.edu/biology/plant/people |
| 6 | www.stanford.edu/chemistry |

- To each URL, a unique integer identifier is assigned.

# Compression of Adjacency Table

| | |
|---|---|
| **nate.com/a.htm** | nate.com/w.htm, naver.com/x.htm |
| **nate.com/b.htm** | nate.com/cafe/boss.htm, nate.com/y.htm |
| **nate.com/c.htm** | nate.com/cafe/boss.htm, naver.com/x.htm |
| **nate.com/d.htm** | nate.com/w.htm, nate.com/y.htm, ever.com/z.htm |

| | |
|---|---|
| **1** | nate.com/a.htm |
| **2** | nate.com/b.htm |
| **3** | nate.com/c.htm |
| **4** | nate.com/d.htm |
| **5** | nate.com/cafe/boss.htm |
| **6** | nate.com/w.htm |
| **7** | naver.com/x.htm |
| **8** | nate.com/y.htm |
| **9** | ever.com/z.htm |

| | |
|---|---|
| **1** | 6, 7 |
| **2** | 5, 8 |
| **3** | 5, 7 |
| **4** | 6, 8, 9 |

# Compression of Adjacency Table

- Main idea : **Similarity**
  The adjacency list of a node is <u>similar</u> to one of the
  <u>7 preceding URLs</u> in the lexicographic ordering.

- Consider the following adjacency table.

| 1 | 1, 2, 4, 8, 16, 32, 64 |
|---|---|
| 2 | 1, 4, 9, 16, 25, 36, 49, 64 |
| 3 | 1, 2, 3, 5, 8, 13, 21, 34, 55, 89, 144 |
| 4 | 1, 4, 8, 16, 25, 36, 49, 64 |

# Compression of Adjacency Table

- Express adjacency list in terms of one of the
  <u>7 preceding URLs</u>.

| | |
|---|---|
| **1** | 1, 2, 4, 8, 16, 32, 64 |
| **2** | **1, 4, 9, 16, 25, 36, 49, 64** |
| **3** | 1, 2, 3, 5, 8, 13, 21, 34, 55, 89, 144 |
| **4** | **1, 4, 8, 16, 25, 36, 49, 64** |

offset -2, remove 9, add 8

2          9      8

# Compression of Adjacency Table

- Main Idea : **Locality**
  Use gap encodings in sorted lists.

| 1 | 1, 2, 4, 8, 16, 32, 64 |
|---|---|
| 2 | **1, 4, 9, 16, 25, 36, 49, 64** |
| 3 | 1, 2, 3, 5, 8, 13, 21, 34, 55, 89, 144 |
| 4 | **1, 4, 8, 16, 25, 36, 49, 64** |

| 1 | 1, +1, +2, +4, +8, +16, +32 |
|---|---|
| 2 | 1, +3, +5, +7, +9, +11, +13, +15 |
| 3 | 1, +1, +1, +2, +3, +5, +8, +13, +21, +34, +55 |
| 4 | offset -2, remove 9, add 8 |

# Compression of Adjacency Table

- How to process connectivity queries?

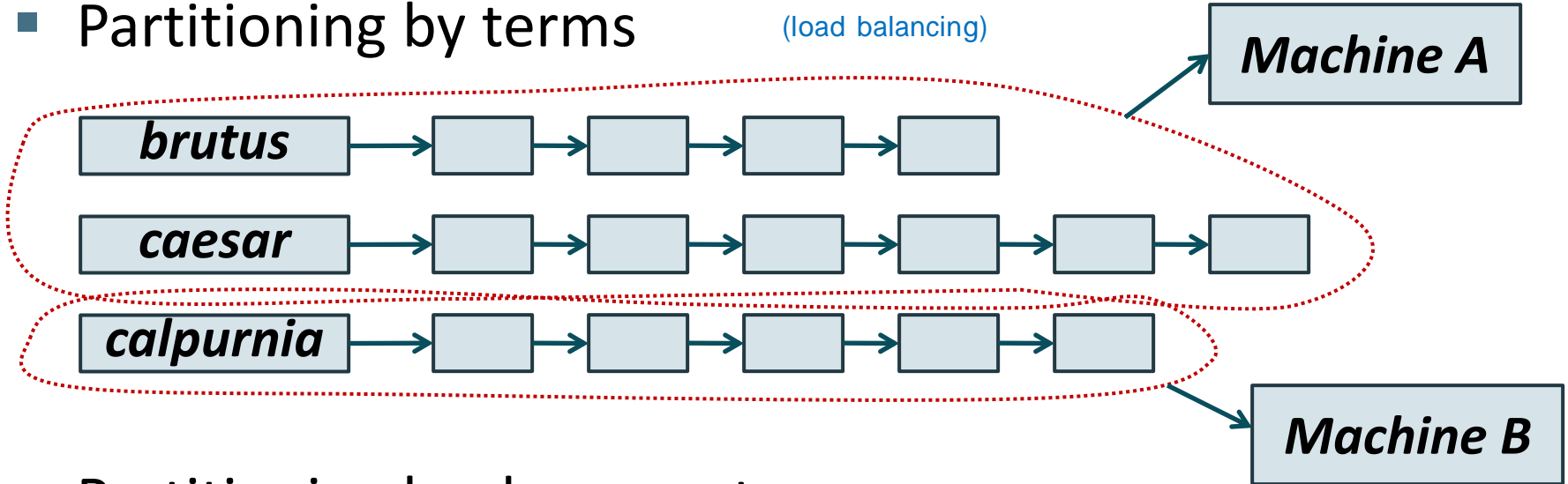| 1 | 1, +1, +2, +4, +8, +16, +32 |
|---|---|
| 2 | 1, +3, +5, +7, +9, +11, +13, +15 |
| 3 | 1, +1, +1, +2, +3, +5, +8, +13, +21, +34, +55 |
| 4 | offset -2, remove 9, add 8 |

- We need to reconstruct the entries of the adjacency table.
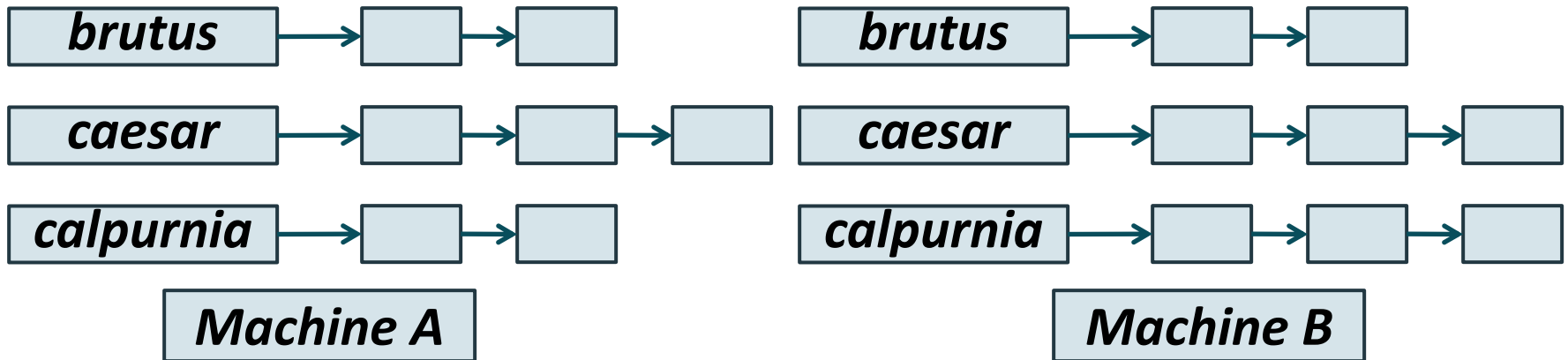
# Distributing Indexes
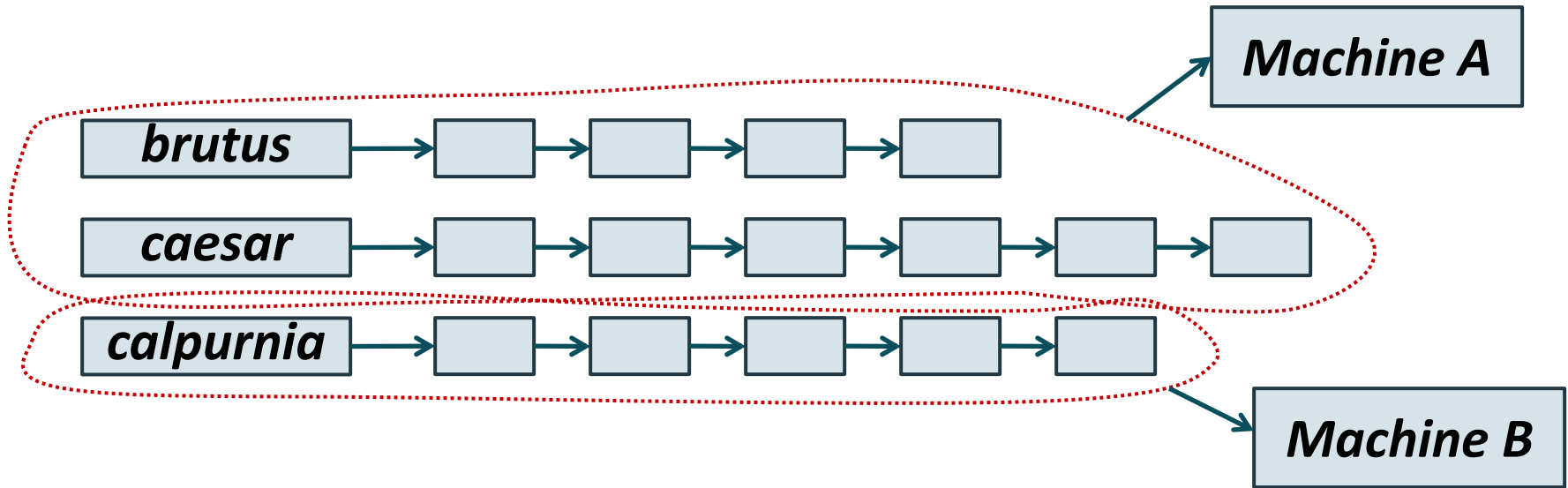
# Distributing Indexes

- Partitioning by terms    (load balancing)

**brutus** → □ → □ → □ → □

**caesar** → □ → □ → □ → □ → □ → □

→ **Machine A**

**calpurnia** → □ → □ → □ → □ → □

→ **Machine B**

- Partitioning by documents

**brutus** → □ → □

**caesar** → □ → □ → □

**calpurnia** → □ → □

**Machine A**

**brutus** → □ → □

**caesar** → □ → □ → □

**calpurnia** → □ → □ → □

**Machine B**

# Partitioning by terms
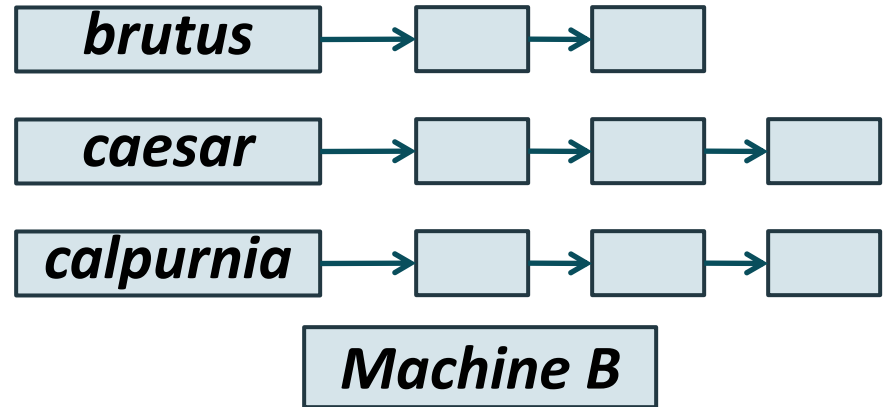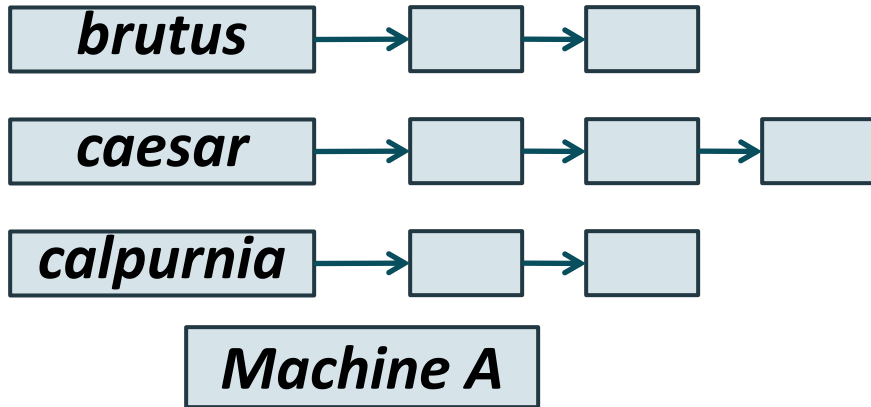
# Partitioning by terms

- The dictionary of index terms is partitioned into subsets.
- A query is routed to the nodes corresponding to its query terms.
- **Pros**
  - Concurrent query processing is possible since different query terms would hit different set of machines.
- **Cons**
  - Multi-word query processing is inefficient due to the overhead of sending long postings list between machines for merging.
  - Load balancing is difficult since it depends on the distribution of query terms not on term frequency.

    load

# Partitioning by documents

# Partitioning by documents

- Each node contains the index for a subset of all documents.
    - Each query is distributed to all nodes.
    - The results from various nodes are merged.

- Multi-word query processing is rather efficient since the merged lists are moving between machines for final merging.

# Partitioning by documents

- How to partition documents to nodes?
    - According to crawling architecture, one simple approach would be to assign *all pages from a **host*** to a single **node**.
        - At query time, the top $k$ results from each **node** being merged to find the top $k$ document for the query.
        - The problem is that search result would come from a small number of **hosts**.
    - A hash of each URL into the space of nodes results in a more **uniform distribution** of query-time computation across nodes.

# Distributing the crawler



document fingerprints

Robot templates

URL set

DNS

www http

Parse

link

text

Content Seen?

URL Filter

Host splitter

To other nodes

hash is used.

From other nodes

Duplicate URL eliminate

Fetch

Indexer

URL Frontier