CS3245

# Information Retrieval

Lecture 5: Index Construction
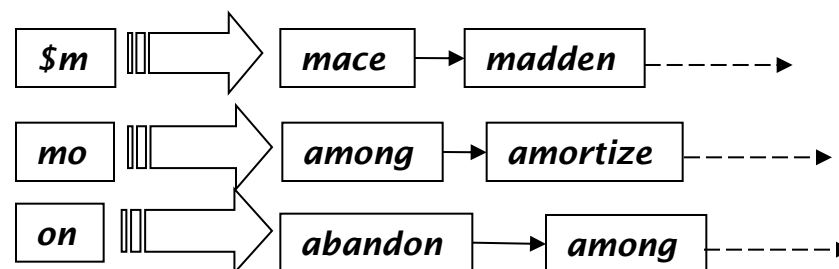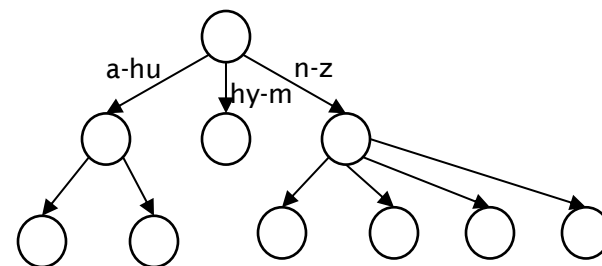
5

# Last Time

- Dictionary data structures

- Tolerant retrieval
  - Wildcards
  - Spelling correction

# Today: Index construction

- How to make index construction scalable?
  1. BSBI  (simple method)
  2. SPIMI (more realistic)
  3. Distributed Indexing

- How to handle changes to the index?
  1. Dynamic Indexing

# Hardware basics

Many design decisions in information retrieval are based on the characteristics of hardware

Especially with respect to the bottleneck:

# **Storage**

# Hardware basics

- Memory is extremely **fast** but **limited in quantity**.
  - DDR5: ~50**GB**/s
  - Available at the magnitude of **GB**

- Hard disk space is **abundant** but a lot **slower**
  - Available at the magnitude of **TB**
  - HDD: ~150 **MB**/s
  - SSD SATA: ~550**MB**/s, SSD NVMe: ~5**GB**/s

# Hardware basics

- Hard disk operations
    - Seek – to access a random location **(no data transfer)**
    - Transfer – to transfer a data block
    - Seek Time >> Transfer Time in general

- Important consideration
    - Data access with a lot of seeks is a huge waste of time!
    - Better to transfer one large (sequential) chunk of data than many small (scattered) chunks.

# Hardware assumptions (Hard Disk)

| symbol | statistic | value |
|--------|-----------|-------|
| s | average seek time | 8 ms = $8 \times 10^{-3}$ s |
| b | transfer time per byte | 0.006 μs = $6 \times 10^{-9}$ s |

# Hardware assumptions (Flash SSDs)

| symbol | statistic | value |
|--------|-----------|-------|
| s | average seek time | .1 ms = $1 \times 10^{-4}$ s |
| b | transfer time per byte | 0.002 μs = $2 \times 10^{-9}$ s |

100x faster seek,
3x faster transfer time.
(But price 8x more per GB of storage)

# RCV1: Our collection for this lecture

- The successor to the Reuters-21578, which you used for your homework assignment. Larger by 35 times.

  - Not really large enough either, but it is publicly available and is a more plausible example.

- One year of Reuters newswire (part of 1995 and 1996)

# Reuters RCV1 statistics

| symbol | statistic | value |
|--------|-----------|-------|
| N | documents | 800,000 |
| L | avg. # tokens per doc | 200 |
| M | terms | 400,000 |
|   | (= vocabulary size) | |
|   | avg. # bytes per term | 7.5 |
| T | term-docID pairs | 100,000,000 |
|   | (= tokens) | |

# Recap: Index Construction (Week 2)

- ## Sort by terms
  - ### And then docID

We focus on this sort step.
We have 100M pairs to sort.

| Term | docID |
|------|-------|
| I | 1 |
| did | 1 |
| enact | 1 |
| julius | 1 |
| caesar | 1 |
| I | 1 |
| was | 1 |
| killed | 1 |
| i' | 1 |
| the | 1 |
| capitol | 1 |
| brutus | 1 |
| killed | 1 |
| me | 1 |
| so | 2 |
| let | 2 |
| it | 2 |
| be | 2 |
| with | 2 |
| caesar | 2 |
| the | 2 |
| noble | 2 |
| brutus | 2 |
| hath | 2 |
| told | 2 |
| you | 2 |
| caesar | 2 |
| was | 2 |
| ambitious | 2 |

| Term | docID |
|------|-------|
| ambitious | 2 |
| be | 2 |
| brutus | 1 |
| brutus | 2 |
| capitol | 1 |
| caesar | 1 |
| caesar | 2 |
| caesar | 2 |
| did | 1 |
| enact | 1 |
| hath | 1 |
| I | 1 |
| I | 1 |
| i' | 1 |
| it | 2 |
| julius | 1 |
| killed | 1 |
| killed | 1 |
| let | 2 |
| me | 1 |
| noble | 2 |
| so | 2 |
| the | 1 |
| the | 2 |
| told | 2 |
| you | 2 |
| was | 1 |
| was | 2 |
| with | 2 |

# Scaling index construction

- At **~11.5** bytes per pair: ~7.5 bytes for term + 4 bytes for docID

- T = 100,000,000 in the case of RCV1: ~1.1GB
  - So ... we can do this easily in memory nowadays, but typical collections are much larger.  E.g. the *New York Times* provides an index of >150 years of newswire

- Thus, we need to make use of the **harddisk**.

# BSBI: Blocked sort-based Indexing

- 8-byte (4+4) records *(termID, docID).*

  - 400,000 terms

  - Create a **dictionary** to map
    **terms** to **termIDs** of 4 bytes

| Term | TermID |
|------|--------|
| noble | 22 |
| Caesar | 1250 |
| killed | 952 |
| Brutus | 3391 |
| … | … |

- These are generated as we parse docs.


- Must now sort 100M 8-byte records by *termID*.

# BSBI: Blocked sort-based Indexing

- Define a <u>Block</u> as ~ **10M** such records
  - Can easily fit a couple into memory.
  - Will have **10** such blocks for our collection.

- Basic idea of algorithm:
  - Map the terms to term ID during generation
  - Accumulate records for each block, **sort**, create the posting lists and write to disk.
  - **Merge** the blocks into bigger blocks recursively.
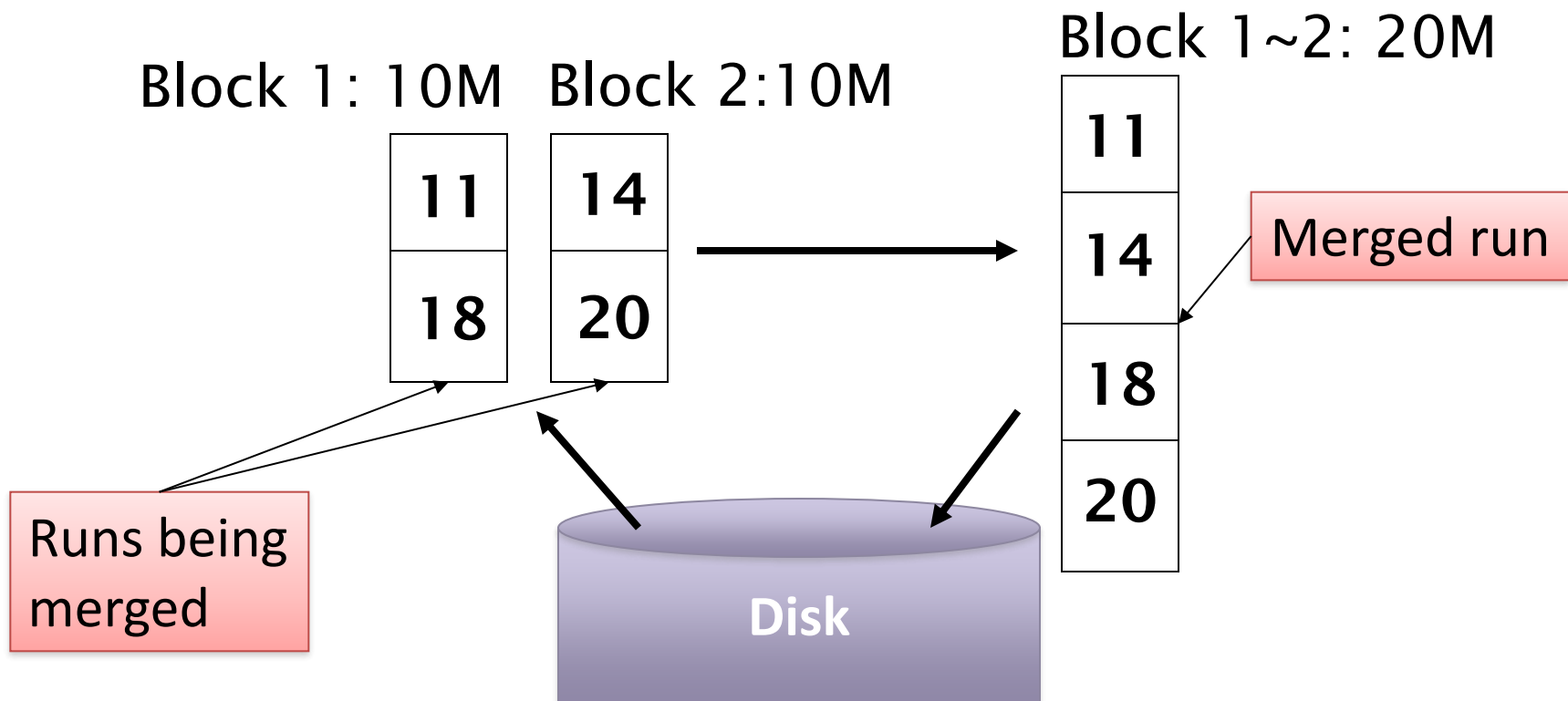  - Reverse the mapping

BSBIndexConstruction()
1    $n \leftarrow 0$
2    **while**   (all documents have not been processed)
3    **do** $n \leftarrow n + 1$
4        $block \leftarrow$ ParseNextBlock()
5        BSBI-Invert($block$)
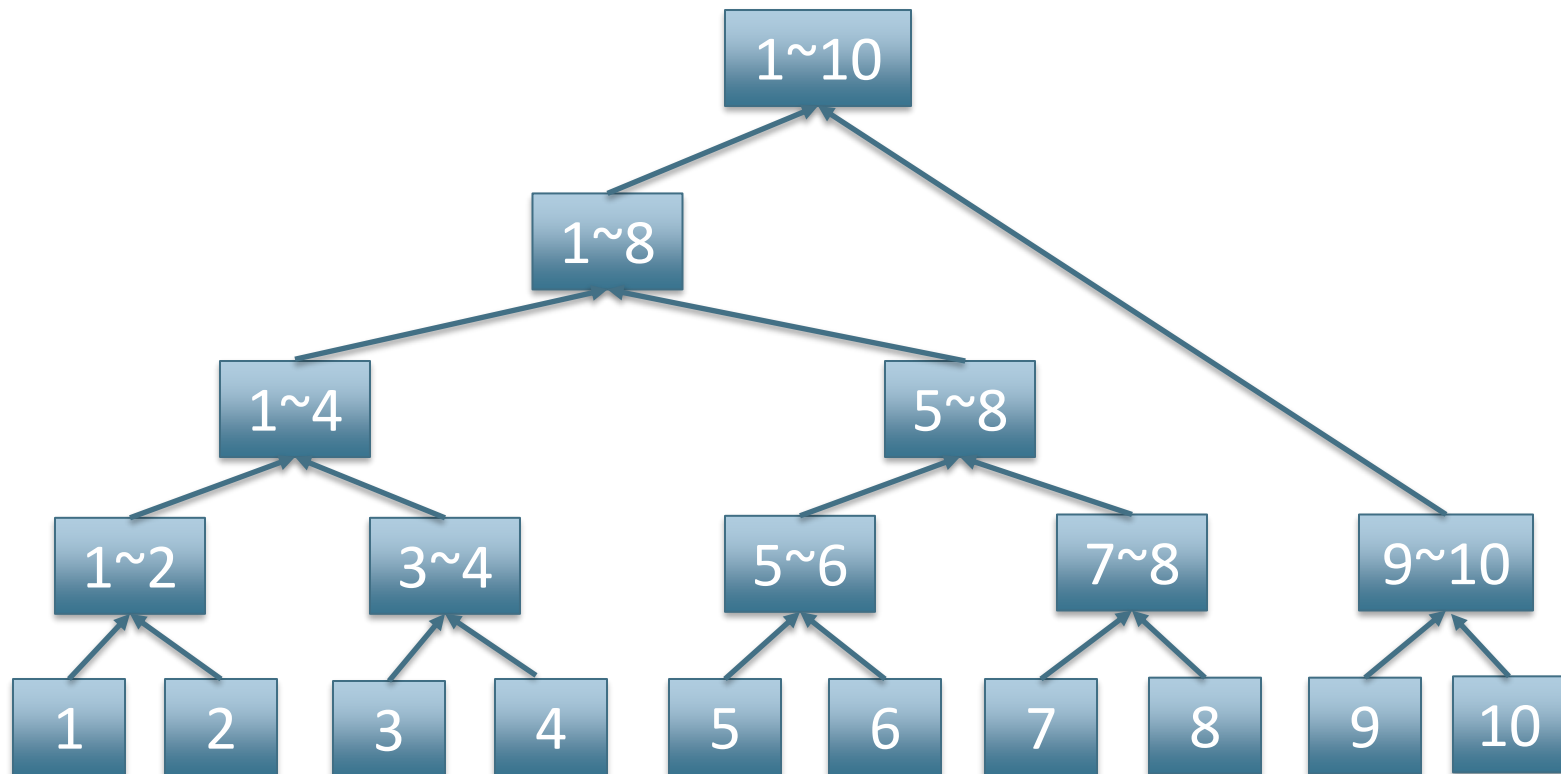6        WriteBlockToDisk($block, f_n$)
7    MergeBlocks($f_1, \ldots, f_n; f_{\text{merged}}$)

# How to merge the sorted runs?

- Can do binary merges,

- Read into memory runs **in blocks of 10M**, merge, write back**.**

Block 1: 10M   Block 2:10M   Block 1~2: 20M

| Block 1 |
|---|
| 11 |
| 18 |

| Block 2 |
|---|
| 14 |
| 20 |

| Merged |
|---|
| 11 |
| 14 |
| 18 |
| 20 |

Merged run

Runs being merged

**Disk**
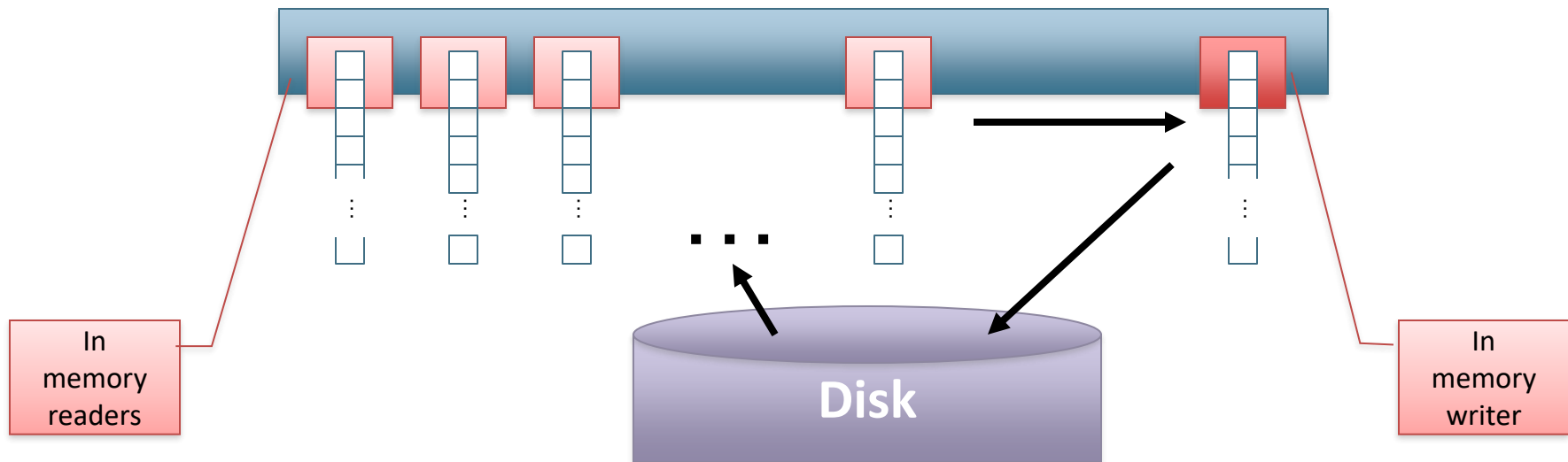
# How to merge the sorted runs?

- 2-way Merge: Merge tree of $\log_2 10 \approx 4$ layers.
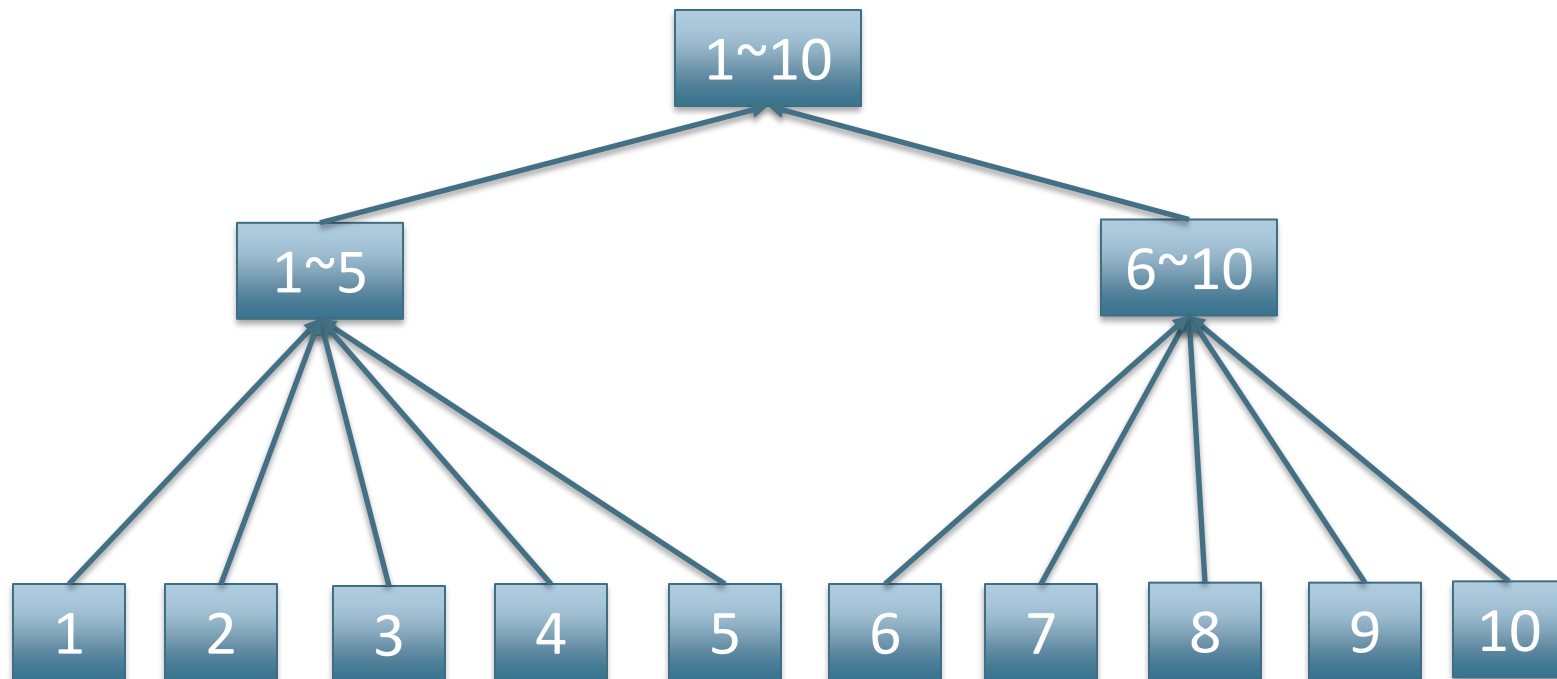
# How to merge the sorted runs?

Second method (better):

▪ It is more efficient to do a *n*-way merge, where you are reading from all blocks simultaneously

▪ Providing you read **decent-sized chunks** of each block into memory and then write out **a decent-sized output chunk**, then your efficiency isn't lost by disk seeks



In memory readers

Disk

In memory writer

# How to merge the sorted runs?

- 5-way Merge: Merge tree of $\log_5 10$ = 2 layers.

# Remaining problems with BSBI

- The dictionary must fit into memory
  - Hard to guarantee since it grows dynamically
  - May end up crashing if the dictionary is too big

- A fixed block size must be decided in advance
  - Too small: could be slow since more blocks need to be processed.
  - Too big: may end up crashing if too much memory is used by other applications.

# SPIMI:
# Single-pass in-memory indexing

- **Key idea 1**: Generate an index (i.e., a **real** dictionary + postings lists) as the pairs are processed

- **Key idea 2**: Go as far as memory allows, write out the index and then merge later

- Advantages:
  - No need to keep a single dictionary in memory
  - No need to wait for a fixed-size block to be filled up
  - Able to adapt to the availability of memory

# SPIMI:
# Single-pass in-memory indexing

| Term | Doc # |
|------|-------|
| I | 1 |
| did | 1 |
| enact | 1 |
| julius | 1 |
| caesar | 1 |
| I | 1 |
| was | 1 |
| killed | 1 |
| i' | 1 |
| the | 1 |
| capitol | 1 |
| brutus | 1 |
| killed | 1 |
| me | 1 |
| so | 2 |
| let | 2 |
| it | 2 |
| be | 2 |
| with | 2 |
| caesar | 2 |
| the | 2 |
| noble | 2 |
| brutus | 2 |
| hath | 2 |
| told | 2 |
| you | 2 |
| caesar | 2 |
| was | 2 |
| ambitious | 2 |
| … | … |

**Hash** the pairs into the index →

| Term | Postings |
|------|----------|
| be | 2 |
| with | 2 |
| caesar | 1,2 |
| hath | 2 |
| i' | 1 |
| it | 2 |
| enact | 1 |
| julius | 1 |
| killed | 1 |
| let | 2 |
| I | 1 |
| did | 1 |
| brutus | 1,2 |
| capitol | 1 |
| you | 2 |
| me | 1 |
| noble | 2 |
| so | 2 |
| ambitious | 2 |
| the | 1,2 |
| told | 2 |
| was | 1,2 |

Inverted Index
(Hash Table
in memory)

**Sort terms** and write out →

| Term | Postings |
|------|----------|
| ambitious | 2 |
| be | 2 |
| brutus | 1,2 |
| caesar | 1,2 |
| capitol | 1 |
| did | 1 |
| enact | 1 |
| hath | 2 |
| I | 1 |
| i' | 1 |
| it | 2 |
| julius | 1 |
| killed | 1 |
| let | 2 |
| me | 1 |
| noble | 2 |
| so | 2 |
| the | 1,2 |
| told | 2 |
| was | 1,2 |
| with | 2 |
| you | 2 |

Inverted Index
(Files on disk)

Merge with other blocks later →

# SPIMI-Invert

SPIMI-INVERT(*token_stream*)
1    *output_file* = NEWFILE()
2    *dictionary* = NEWHASH()
3    **while**  (free memory available)
4    **do** *token* ← *next*(*token_stream*)
5        **if** *term*(*token*) ∉ *dictionary*
6            **then** *postings_list* = ADDTODICTIONARY(*dictionary*, *term*(*token*))
7            **else**  *postings_list* = GETPOSTINGSLIST(*dictionary*, *term*(*token*))
8        **if** *full*(*postings_list*)
9            **then** *postings_list* = DOUBLEPOSTINGSLIST(*dictionary*, *term*(*token*))
10        ADDTOPOSTINGSLIST(*postings_list*, *docID*(*token*))
11    *sorted_terms* ← SORTTERMS(*dictionary*)
12    WRITEBLOCKTODISK(*sorted_terms*, *dictionary*, *output_file*)
13    **return** *output_file*

- Merging of blocks is analogous to BSBI.

# SPIMI: Efficiency

- Faster than BSBI
    - No sorting of pairs
    - Only sorting of dictionary terms

- Even faster with compression
    - Compression of terms
    - Compression of postings

More about this in W6.

# DISTRIBUTED INDEXING

# Distributed indexing

- For web-scale indexing (don't try this at home!):
    must use a distributed computing cluster

- Individual machines are fault-prone
    Can unpredictably slow down or fail

How do we exploit such a pool of machines?

# Google Data Centers

- Google data centers mainly contain commodity machines, and are distributed worldwide.

- One here in Jurong West (~200K servers back in 2011)

- Must be fault tolerant.  Even with 99.9+% uptime, there often will be one or more machines down in a data center.

- As of 2001, they have fit their entire web index in-memory (RAM; of course, spread over many machines)
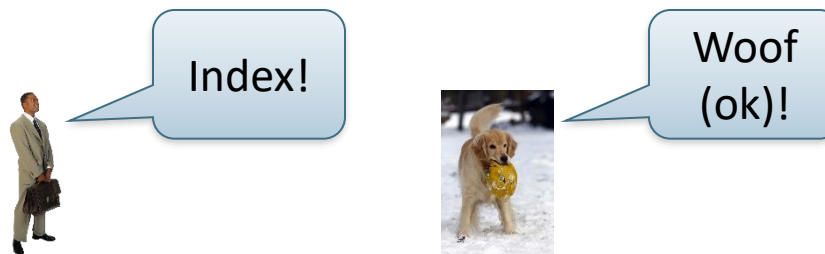


https://youtu.be/XZmGGAbHqa0

http://www.gizmodo.com.au/2010/04/googles-insane-number-of-servers-visualised/

http://www.google.com/about/datacenters/inside/streetview/

http://www.straitstimes.com/business/10-things-you-should-know-about-google-data-centre-in-jurong

# Architecture of distributed indexing

- Maintain a *master* machine directing the indexing job – considered "safe".
  - Master nodes can fail too!
- Break up indexing into sets of (parallel) tasks.
- Master machine assigns each task to an idle *worker* machine from a pool.

Index!

Woof (ok)!

# Parallel tasks

- We will use two sets of parallel tasks
  - Parsers    Parse documents and emit pairs
  - Inverters    Sort the pairs and build the index

- Break the input document collection into *splits*
- Each split is a subset of documents (corresponding to blocks in BSBI/SPIMI)

# Parsers

- Master assigns a split to an idle parser machine

- Parser reads a document at a time and emits (term, doc) pairs

- Parser writes pairs into $j$ partitions

- Each partition is for a range of terms' first letters
  - (e.g., *a-f, g-p, q-z*) – here $j$ = 3.
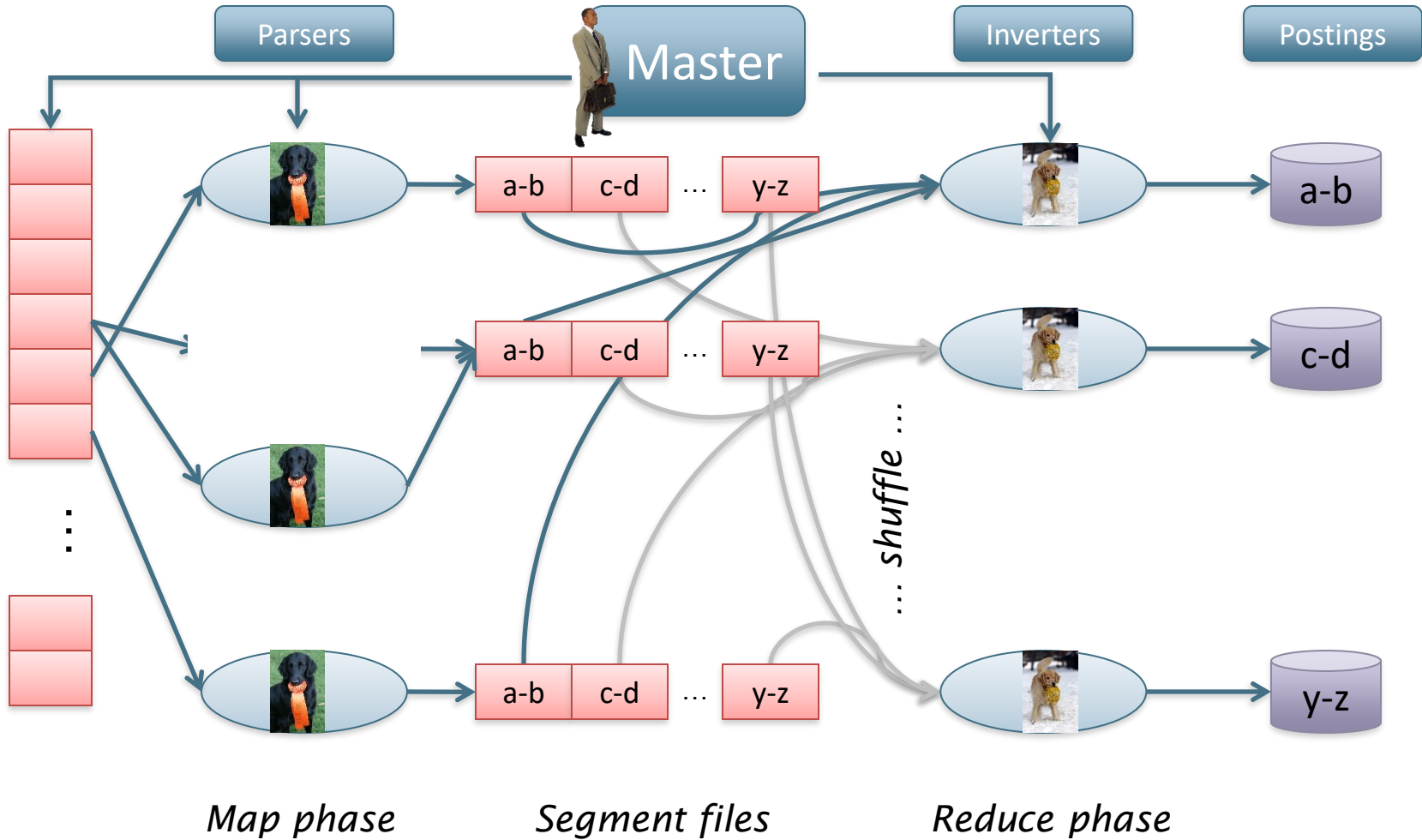  - (e.g., *a-b, c-d, …, y-z*) – here $j$ = 13.

# Inverters

- Master assigns a term-partition to an idle inverter machine

- Inverter collects all (term,doc) pairs (= postings) from the partition.

- Inverter sorts and writes to postings lists

# Data flow



Parsers　　Master　　Inverters　　Postings

| a-b | c-d | … | y-z |

| a-b | c-d | … | y-z |

| a-b | c-d | … | y-z |

a-b

c-d

y-z

*… shuffle …*

*Map phase*　　　*Segment files*　　　*Reduce phase*

# MapReduce

- The index construction algorithm we just described is an instance of MapReduce.

- MapReduce is a robust and conceptually simple framework for distributed computing
  … without having to write code for the distribution part.

- They describe the Google indexing system as consisting of a number of phases, each implemented in MapReduce.

# MapReduce for indexing

**Schema of map and reduce functions**

- **map**: input → list(k, v)    **reduce**: (k, list(v)) → output

**Instantiation of the schema for index construction**

- **map**: **blocks of** web collection → list(term, docID) **in segment files**

- **reduce**: (<term1, list(docID)>, <term2, list(docID)>, …) **from segment files for the same partition → consolidated blocks of** (term1: postings list1, term 2: postings list2, …)

# MapReduce

- **map**

  - d1 : Caesar came, Caesar conquered. d2 : Caesar died →

  - (caesar, d2), (died,d2), (caesar, d1), (came, d1), (caesar, d1), (conquered, d1)

- **Reduce**

  - <caesar, (d2, d1, d1)>, <died, (d2)>, <came, (d1)>, <conquered, (d1)> →

  - <caesar, (d1, d2)>, <came, (d1)>, <conquered, (d1)>, <died, (d2)>

# DYNAMIC INDEXING

# Dynamic indexing

- Up to now, we assumed that collections are static.

- In practice, they rarely are!
    - Documents come in over time and need to be inserted.
    - Documents are deleted and modified.

- The dictionary and postings lists have to be modified:
    - Postings updates for terms already in dictionary
    - New terms added to dictionary

- First approach: re-index every time?
    - Simple but impractical

# 2$^{nd}$ simplest approach

- Two indexes
  - One "big" main index (let say **I**)
  - One "small" (in memory) auxiliary index (let say **Z**)

- Mechanism
  - Add: new docs goes to the auxiliary index
  - Delete: maintain a list of deleted docs
  - Update: delete + add
  - Search: search both, merge results and omit deleted docs

- Need to perform **linear merge** when auxiliary index is too large.

# Linear Merge

- ## Let say…

    - The capacity of the auxiliary index **Z** is **n** pairs of (term, docID)

    - The main index **I** can be arbitrarily large

    - Initially both are empty

- ## The algorithm

    - Once Z is full, write out **Z** and merge with **I**

# Linear Merge

- Example:
    - The 1$^{st}$ set of **n** pairs, write out **Z** (**n** items) and merge with **I** (**0** items) → merge **n + 0 = n** items into **I**
    - The 2$^{nd}$ set of **n** pairs, write out **Z** (**n** items) and merge with **I** (**n** items) → merge **n + n** = **2*n** items into **I**
    - The 3$^{rd}$ set of **n** pairs, write out **Z** (**n** items) and merge with **I** (**2*n** items) → merge **n + 2*n = 3*n** items into **I**
    - The 4$^{th}$ set of **n** pairs, write out **Z** (**n** items) and merge with **I** (**3*n** items) → merge **n + 3*n = 4*n** items into **I**
    - …

# Linear Merge

- Let say there are a total **T** pairs for which require **k** merges (i.e., k = T / n)

- Cost of merging
  - n + 2 * n + 3 * n + 4 * n ... + k * n
    = (k * (k+1) / 2) * n
    ~= $nk^2$
    ~= **$O(T^2)$**

# Logarithmic merge

- Idea: maintain a series of indexes
  - Z: In memory, with the same capacity as $I_0$ (= $n$)
  - $I_0$, $I_1$, …: on disk, each twice as large as the previous one.

*Loop for log levels*

  - If Z gets too big (= $n$), write to disk as $I_0$,
    or merge with $I_0$ (if $I_0$ already exists) as $I_1$
  - Either write $I_1$ to disk as $I_1$,
    or merge with $I_1$ (if $I_1$ already exists) to form $I_2$
    … etc.

# Logarithmic merge

- Example:
  - The 1$^{st}$ set of **n** pairs, write out **Z** (**n** items) as **I$_0$**
  - The 2$^{nd}$ set of **n** pairs, write out **Z** (**n** items) but **I$_0$** already exists → merge **n + n** = **2\*n** items into **I$_1$** (and **I$_0$** is gone)
  - The 3$^{rd}$ set of **n** pairs, write out **Z** (**n** items) as **I$_0$**
  - …

|       | I$_0$ | I$_1$ | I$_2$ |
|-------|-------|-------|-------|
| **0**   | 0 | 0 | 0 |
| **n**   | 1 | 0 | 0 |
| **2\*n** | 0 | 1 | 0 |
| **3\*n** | 1 | 1 | 0 |
| **4\*n** | 0 | 0 | 1 |

The presence (1) or absence (0) of the indexes on disk

# Logarithmic merge

- Example:
  - …
  - The 4th set of **n** pairs, write out **Z** (**n** items) but **I$_0$** already exists → merge **n + n** = **2\*n** items into a new index **I$_1$** but **I$_1$** already exists → merge **2\*n + 2\*n = 4\*n** items into a new index **I$_2$** (and **I$_0$** and **I$_1$** are gone).

|       | I$_0$ | I$_1$ | I$_2$ |
|-------|-------|-------|-------|
| **0**   | 0 | 0 | 0 |
| **n**   | 1 | 0 | 0 |
| **2\*n** | 0 | 1 | 0 |
| **3\*n** | 1 | 1 | 0 |
| **4\*n** | 0 | 0 | 1 |

The presence (1) or absence (0) of the indexes on disk

LMERGEADDTOKEN(*indexes*, $Z_0$, *token*)
1　$Z_0 \leftarrow$ MERGE($Z_0$, {*token*})
2　**if** $|Z_0| = n$
3　　**then for** $i \leftarrow 0$ **to** $\infty$
4　　　　**do if** $I_i \in$ *indexes*
5　　　　　　**then** $Z_{i+1} \leftarrow$ MERGE($I_i$, $Z_i$)
6　　　　　　　　(*$Z_{i+1}$ is a temporary index on disk.*)
7　　　　　　　　*indexes* $\leftarrow$ *indexes* $- \{I_i\}$
8　　　　　**else** $I_i \leftarrow Z_i$　　(*$Z_i$ becomes the permanent index $I_i$.*)
9　　　　　　　*indexes* $\leftarrow$ *indexes* $\cup \{I_i\}$
10　　　　　　　BREAK
11　　　　$Z_0 \leftarrow \emptyset$


LOGARITHMICMERGE()
1　$Z_0 \leftarrow \emptyset$　　(*$Z_0$ is the in-memory index.*)
2　*indexes* $\leftarrow \emptyset$
3　**while** true
4　**do** LMERGEADDTOKEN(*indexes*, $Z_0$, GETNEXTTOKEN())

# Logarithmic merge

- ## Cost of merging

  - Each posting is touched *O(log* T) times, so complexity is *O(T log T)*

  - E.g., let n = 4, T = 32, the first pair is touched 4 times (as compared to 8 times in linear merge)

- ## So logarithmic merge is much more efficient for indexing

- ## But query processing now is slower

  - Merging results from *O(log T)* indexes (as compared to 2)

# Summary

- Indexing
  - Both basic as well as important variants
    - BSBI – sort key values to merge, needs dictionary
    - SPIMI – build mini indexes and merge them, no dictionary
  - Distributed
    - Described MapReduce architecture – a good illustration of distributed computing
  - Dynamic
    - Tradeoff between querying and indexing complexity

# Resources for today's lecture

- Chapter 4 of IIR

- MG Chapter 5

- Original publication on MapReduce: Dean and Ghemawat (2004)

- Original publication on SPIMI: Heinz and Zobel (2003)